



HAL
open science

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

Bagdevi Mishra, Bartosz Ulaszewski, Joanna Meger, Jean-Marc Aury, Catherine Bodénès, Isabelle Lesur, Markus Pfenninger, Corinne da Silva, Deepak Gupta, Erwan Guichoux, et al.

► To cite this version:

Bagdevi Mishra, Bartosz Ulaszewski, Joanna Meger, Jean-Marc Aury, Catherine Bodénès, et al.. A Chromosome-level genome assembly of the European Beech (*Fagus sylvatica*) reveals anomalies for organelle DNA integration, repeat Content and distribution of SNPs. *Frontiers in Genetics*, 2022, 12, pp.691058. 10.3389/fgene.2021.691058 . hal-03602006

HAL Id: hal-03602006

<https://hal.inrae.fr/hal-03602006>

Submitted on 22 Jun 2022

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License



A Chromosome-Level Genome Assembly of the European Beech (*Fagus sylvatica*) Reveals Anomalies for Organelle DNA Integration, Repeat Content and Distribution of SNPs

Bagdevi Mishra^{1,2}, Bartosz Ulaszewski³, Joanna Meger³, Jean-Marc Aury⁴, Catherine Bodénès⁵, Isabelle Lesur-Kupin^{5,6,7}, Markus Pfenninger¹, Corinne Da Silva⁴, Deepak K Gupta^{1,2,8}, Erwan Guichoux⁵, Katrin Heer^{7,9}, Céline Lalanne⁵, Karine Labadie⁴, Lars Opgenoorth⁷, Sebastian Ploch¹, Grégoire Le Provost⁵, Jérôme Salse¹⁰, Ivan Scotti¹¹, Stefan Wötzel^{1,2}, Christophe Plomion⁵, Jaroslaw Burczyk³ and Marco Thines^{1,2,8*}

OPEN ACCESS

Edited by:

Frank M. You,
Agriculture and Agri-Food Canada
(AAFC), Canada

Reviewed by:

Daniel G. Peterson,
Mississippi State University,
United States
Pingchuan Li,
Agriculture and Agri-Food Canada
(AAFC), Canada

*Correspondence:

Marco Thines
marco.thines@senckenberg.de

Specialty section:

This article was submitted to
Plant Genomics,
a section of the journal
Frontiers in Genetics

Received: 05 April 2021

Accepted: 14 December 2021

Published: 08 February 2022

Citation:

Mishra B, Ulaszewski B, Meger J, Aury J-M, Bodénès C, Lesur-Kupin I, Pfenninger M, Da Silva C, Gupta DK, Guichoux E, Heer K, Lalanne C, Labadie K, Opgenoorth L, Ploch S, Le Provost G, Salse J, Scotti I, Wötzel S, Plomion C, Burczyk J and Thines M (2022) A Chromosome-Level Genome Assembly of the European Beech (*Fagus sylvatica*) Reveals Anomalies for Organelle DNA Integration, Repeat Content and Distribution of SNPs. *Front. Genet.* 12:691058. doi: 10.3389/fgene.2021.691058

¹Senckenberg Biodiversity and Climate Research Centre (BiK-F), Senckenberg Gesellschaft für Naturforschung, Frankfurt am Main, Germany, ²Department for Biological Sciences, Institute of Ecology, Evolution and Diversity, Goethe University, Frankfurt am Main, Germany, ³Department of Genetics, ul. Chodkiewicza 30, Kazimierz Wielki University, Bydgoszcz, Poland, ⁴Genoscope, Institut François Jacob, CEA, CNRS, Univ Evry, Université Paris-Saclay, Evry, France, ⁵INRAE, Univ. Bordeaux, BIOGECO, Cestas, France, ⁶HelixVenture, Mérignac, France, ⁷Faculty of Biology, Plant Ecology and Geobotany, Philipps University Marburg, Marburg, Germany, ⁸LOEWE Centre for Translational Biodiversity Genomics (TBG), Frankfurt am Main, Germany, ⁹Forest Genetics, Albert-Ludwigs-Universität Freiburg, Freiburg, Germany, ¹⁰INRAE, URFM, Avignon, France, ¹¹INRAE, UCA, GDEC, Clermont-Ferrand, France

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 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 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.

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

1 INTRODUCTION

Many lowland and mountainous forests in Central Europe are dominated by the European Beech (*Fagus sylvatica*) (Durrant et al., 2016). European Beech 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 et al., 2006; Geßler et 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 et al., 2015), will intensify climatic stress as already now severe damage has been observed (Geßler et al., 2007; Reif et al., 2017). In order to cope with this, and 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 life, including animals (Michael and VanBuren, 2020; Priest et al., 2020; Rhie et 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 et 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 et al., 2012). For this family, various draft assemblies have been published (Martínez-García et al., 2016; Plomion et al., 2016; Sork et al., 2016),



FIGURE 1 | The more than 300 year-old *Fagus sylvatica* reference individual Bhaga on a cliff over the Edersee in the Kellerwald Edersee National Park (Germany).

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 (**Figure 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 library construction and sequencing was done by a commercial sequencing provider (BGI, Hong Kong, and 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. (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 reconstruction) and at the National Forest Office nursery in Guémené-Penfao in January 2019 (429 saplings corresponding to the 2nd campaign, used for linkage mapping). For relatedness assessment among the half-sib progeny of MSSB, young leaves after bud burst were sampled from saplings in the nursery in spring 2014 (1st campaign) and 2017 (2nd campaign), immediately frozen in dry ice and then stored at -80°C before subsequent genetic analyses. Likewise, leaves were sampled on the mother tree and 19 surrounding adult trees (expected fathers). Nuclear DNA was extracted individually from 10 mg of tissue using the DNeasy Plant Mini Kit (QIAGEN, Hilden, Germany) following the manufacturer's instructions. DNA concentration was measured on a ND-8000 NanoDrop spectrophotometer (Thermo Scientific, Wilmington, United States). For additional transcriptome construction a total of six different organs were sampled on the MSSB accession, including: two types of buds (quiescent buds and swelling buds just before bud break) during dormancy release on the 15th March, 2017, male flowers and female flowers collected on the 3rd of May 2017, leaves, and xylem collected on the 28th of June 2017. Each organ was immediately flash-frozen in liquid nitrogen and stored at -80°C before RNA extraction. For short read sequencing (Illumina), total RNA was extracted from these six samples following the procedure described in Le Provost et al. (2007). Residual genomic DNA was removed before purification using DNase RQ1 (Promega, Madison, WI, United States) according to the manufacturer's instructions. The quantity and the quality of each extract was determined using an Agilent 2,100 Bioanalyser (Agilent Technologies, Inc., Santa Clara, CA, United States). For long read sequencing (Oxford Nanopore Technologies, Oxford, United Kingdom) total RNA was extracted as described above and depleted using the Ribo-Zero rRNA Removal Kit Plant Leaves (Illumina, San Diego, CA, United States). RNA was then purified and concentrated on a RNA Clean Concentrator™-5 column (Zymo Research, Irvine, CA, United States).

For the linkage mapping, vegetative buds from the individuals from the first and second campaign were sampled on the February 28, 2018 in Nouzilly at the ecodormancy stage from 200 genotypes (i.e., 200 half-sibs that constitute the mapping

population), were frozen on dry ice, and then stored at -80°C. RNA was extracted from bud-scale-free leaves following the procedure described above. These 200 genotypes included two relatively large full-sib families comprising 49 full-sibs (family MSSBxSSP12) and 36 full-sibs (family MSSBxMSSH) (see results section).

2.2 Chromosomal Pseudo-Molecules and Their Annotation

2.2.1 Building of Chromosomal Pseudo-molecules Using Hi-C Reads

The previous scaffold-level assembly was constructed with Illumina shotgun short reads and PacBio long reads (Mishra et al., 2018). For a chromosome-level assembly, intermediate results from the previous assembly were used as the starting material. Sequence homology of the 6,699 scaffolds generated from the DBG2OLC hybrid assembler (Ye et al., 2016), to the separately assembled chloroplast and mitochondria of beech, and were inferred using blast v2.10.1 (Altschul et al., 1990). All scaffolds that match in full length to any of the organelle with identity >99% and gaps and/or mismatches ≤3 were discarded. The remaining 6,657 scaffolds along with Hi-C data were used in ALLHiC (Zhang et al., 2019) for building the initial chromosome-level assembly. The cleaned Illumina reads were aligned to the initial assembly using Bowtie2 software (Langmead and Salzberg, 2012) and then, sorted and indexed bam files of the concordantly aligned read pairs for all the sequences were used in Pilon (Walker et al., 2014) to improve the correctness of the assembly. The final assemblies for Bhaga and Jamy were deposited at NCBI under the accession numbers PRJEB43845 and PRJNA450822, respectively.

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

2.2.2 Gene Prediction

Cleaned transcriptomic Illumina reads (minimum read length: 70; average read quality: 25 and read pairs containing no N) were aligned to the assembly using Hisat (Kim et al., 2015) in order to generate splice-aware alignments. The sorted and indexed bam file [samtools, v1.9 (Li et al., 2009)] of the splice-aware alignments was used in the "Eukaryotic gene finding" pipeline of OmicsBox—Bioinformatics Made Easy (2020) which uses Augustus (Stanke and Morgenstern, 2005) for gene prediction. For predictions, some parameters were changed from the default values, i.e. minimum intron length was set to 20, minimum exon length was set to 200, and complete genes (with start and stop codon) of a minimum of 180 bp length were predicted, by choosing *Arabidopsis thaliana* as the closest reference-level organism.

2.2.3 Assessment of the Gene Space

The protein sequences of the PLAZA (Bell, 2018) genes for *A. thaliana*, *Vitis vinifera*, and *Eucalyptus grandis* were downloaded from the plaza v4.5 dicots dataset (accessed October 21, 2020)

and were used along with the predicted proteins from the current assembly to build protein clusters using cd-hit v.4.8.1 (Li and Godzik, 2006; Fu et al., 2012). The number of exons per genes was assessed and compared to the complete coding genes from *A. thaliana*, *Populus trichocarpa*, and *Castanea mollissima*, in line with the comparison made in the scaffold-level assembly (Mishra et al., 2018).

2.2.4 Functional Annotation of Genes

The predicted genes were translated into proteins using transeq (EMBOSS:6.6.0.0, Rice et al., 2000) and queried against the non-redundant database from NCBI (downloaded on 2020-06-24)¹ using diamond (v0.9.30) (Buchfink and Xie, 2015) to find homology of the predicted proteins to sequences with known functions. For prediction of protein family membership, as well as the presence of functional domains and sites in the predicted proteins, Interproscan v5.39.77 (Jones et al., 2014) was used. Result files from both diamond and Interproscan (in Xml format) were used in the blast2go (Götz et al., 2008) module of OmicsBox, and taking both homology and functional domains into consideration, final functional annotations were assigned to the genes. The density of the coding space for each 100 kb region stretch was calculated for all the chromosomes.

2.2.5 Repeat Prediction and Analysis

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

In a separate analysis, repeat elements present in *Fagus sylvatica* were identified by a combination of homology-based and *de novo* approaches using RepeatModeler 2.0 (Flynn et al., 2018) and RepeatMasker v. 4.1.1 (Tarailo-Graovac and Chen, 2009). First, were repetitive elements were identified and classified *de novo* and generated a library of consensus sequences using RepeatModeler 2.0 (Flynn et al., 2018). Then repeats in the assembly were annotated with RepeatMasker 4.1.1 (Tarailo-Graovac and Chen, 2009) using the custom repeat library generated in the previous step.

2.2.6 Telomeric and Centromeric Repeat Identification

Tandem repeat finder (TRF version 4.0.9) (Benson, 1999) was used with parameters 2, 7, 7, 80, 10, 50, and 500 for Match, Mismatch, Delta, PM, PI, Minscore, and MaxPeriod, respectively (Marrano et al., 2020), and all tandem repeats with a monomer length up to 500 bp were predicted. Repeat frequencies of all monomers were plotted against the length of the monomers to identify high-frequency repeats. As the repeats were fetched by TRF with different start and end positions and, thus, identical

repeats were falsely identified as different ones, the program MARS (Ayad and Pissis, 2017) was used to align the monomers of the different predicted repeats, and the repeat frequencies were adjusted accordingly. The chromosomal locations of telomeric and centromeric repeats were identified by blasting the repeats to the chromosomes. For confirmation of centromeric locations, pericentromeres of *A. thaliana* were blasted against the chromosomes of Bhaga.

2.2.7 Organelle Integration

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

2.2.8 SNP Identification and Assessment

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

2.2.9 Genome Browser

A genome browser was set up using JBrowse v.1.16.10 (Buels et al., 2016). Tracks for the predicted gene models and annotated repeat elements were added using the gff files. Separate tracks for the SNP locations and the locations of telomere and centromere

¹<https://ftp.ncbi.nlm.nih.gov/blast/db/> (Accessed June 24, 2020).

repeats were added as bed files. A track depicting the GC content was also added. The genome browser can be accessed from <http://beechgenome.net>.

2.3 Pedigree Reconstruction

2.3.1 SNP Assay Design and Genotyping for Relatedness Assessment Among Half-Sibs

A multiplexed assay using the MassARRAY® MALDI-TOF platform (iPLEX MassArray, Agena BioScience, United States) was used to genotype the mother tree (MSSB), its half-sib progeny from the 1st campaign, and 19 putative fathers. PCR and extension primers were designed from flanking sequences (60 bp of either side) of 40 loci (**Supplementary File S5**) available from Lalagüe et al. (2014) as well as Ouayjan and Hampe (2018). Data analysis was performed with Typer Analyzer 4.0.26.75 (Agena BioScience). Monomorphic SNPs were filtered out, as well as loci with a weak or ambiguous signal (i.e., displaying more than three clusters of genotypes or unclear cluster delimitation). Thirty-six SNPs were finally retained for the paternity analysis.

2.3.2 Sibship Assignment

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

2.4 Unigene Set Construction

2.4.1 Library Construction and Sequencing

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

One Nanopore cDNA library was also prepared from entire female flower RNA. The cDNA library was obtained from 50 ng RNA following the Oxford Nanopore Technologies (Oxford Nanopore Technologies Ltd., Oxford, United Kingdom) protocol “cDNA-PCR Sequencing (SQK-PCS108)” with a 14 cycle PCR (6 min for elongation time). ONT adapters were ligated to 190 ng of resulting cDNA. The Nanopore library was sequenced using a MinION Mk1b and R9.4.1 flowcells.

2.4.2 Bioinformatic Analysis

Short-read RNA-Seq data (Illumina) from the six tissues were assembled using Velvet (Zerbino and Birney, 2008) 1.2.07 and Oases (Schulz et al., 2012) 0.2.08, using a k-mer size of 63 bp.

Reads were mapped back to the contigs with BWA-mem (Li et al., 2009) and consistent paired-end reads were selected. Chimeric contigs were identified and split (uncovered regions) based on coverage information from consistent paired-end reads. Moreover, open reading frames (ORF) and domains were searched using, respectively, TransDecoder (Haas et al., 2013) and CDDsearch (Marchler-Bauer et al., 2011). Only breaks outside ORF and domains were allowed. Finally, the read strand information was used to correctly orient the RNA-seq contigs.

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

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

2.5 Genotyping-By-Sequencing of the Mapping Population

2.5.1 RNAseq Libraries Construction

The 200 RNA samples were prepared as described above (Unigene Set Construction section), using the TruSeq Stranded mRNA kit (Illumina, San Diego, CA, United States), starting from 500 ng total RNA. Libraries were multiplexed onto an Illumina Novaseq 6,000 sequencer (Illumina, San Diego, CA, United States), using the S4 chemistry (2 × 150 read length), targeting approximately 30 million reads per sample.

2.5.2 RNAseq Read Processing for the MSSB Accession

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

2.5.3 Identification of SNPs From RNAseq Data and Offspring Genotype Inference

SNPs were called and bioinformatically genotyped for the mapping population at each MSSB polymorphic site, based on the paired-end Illumina sequencing of 200 RNAseq libraries. The 200 raw-read datasets were trimmed following the same

procedure used for MSSB. Reads were aligned to the 35 K unigene set using *bwa mem* 0.7.17. Genotypes were recovered from the 200 libraries at the 507,905 polymorphic positions identified in MSSB using GATK 3.8.

Subsequently, the following four-step filtering procedure was applied: 1) for each SNP of a given half-sib, polymorphic genotypes were set to monomorphic if the sequencing depth for this individual at a given position was lower than 20; 2) SNPs were kept only if at least 50% of the mapping population (i.e., 100 half-sibs) were heterozygous at that site; 3) only polymorphic sites consistent with a 1:1 heterozygote:homozygote genotype ratio were kept, according to a Chi-square test with a 90% confidence interval (Chi-square < 6.635, 1 d.f.), corresponding to heterozygous loci in the mother tree and monomorphic in all possible fathers; 4) finally, for each contig, only the SNP with fewest missing data in the mapping population was retained.

2.6 Linkage Map Construction

Half-sibs presenting a high amount of missing data were discarded. As a result, 182 individuals (out of 200 selected from the first and second campaign) with valid genotypes for at least 4,127 loci were kept for further analyses. A preliminary analysis was then performed using R-qtL package to group linked SNP markers into robust linkage groups (LG) (LOD = 8) (**Supplementary File S6**). Given the large number of markers per LG, marker ordering was performed within each LG using JoinMap 4.1 (Kyazma, Wageningen, Netherlands). To this end, linkage groups of the maternal parent (MSSB) were constructed using a four-step procedure: 1) The maximum likelihood (ML) algorithm of JoinMap was first used with a minimum linkage LOD score of 5 to calculate the number of crossing-overs (CO) for each individual and to estimate the position of all mapped SNPs, 2) then, the regression algorithm (with a minimum LOD of 5 and default parameters: recombination frequency of 0.4 and maximum threshold value of 1 for the jump) was used for a subset of evenly spaced SNPs (referred to below as set #1 SNPs) along each LG, 3) the maternal linkage maps of the two full-sib families, identified from the paternity test, were constructed using this subset of markers and individuals, providing two genetic maps (referred to below as set #2 and set #3 SNPs) with higher confidence in genetic distance estimates and marker ordering, and both parents being known; 4) finally, from these two SNP datasets, a final dataset was created (set #4) combining sets #2 and #3. For these 3 marker sets (#2, #3, and #4), a first map was constructed using the ML algorithm to calculate the number of CO, and a second map was established using the regression algorithm excluding SNPs with high conflict of positions, reducing the number of CO.

2.7 Genomic Scaffold Anchoring

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

3 RESULTS

3.1 General Genome Features

3.1.1 Genomic Composition and Completeness

The final assembly of the Bhaga genome was based on hybrid assembly of PacBio and Illumina reads as well as scaffolding using a Hi-C library. It was resolved into 12 chromosomes, spanning 535.4 Mb of the genome and 155 unassigned contigs of 4.9 Mb that to 79% consisted of unplaced repeat regions that precluded their unequivocal placement. The genome revealed a high level of BUSCO gene detection (97.4%), surpassing that of the previous assembly and other genome assemblies available for members of the *Fagaceae* (**Table 1**). Of the complete genome assembly, 57.12% were annotated as interspersed repeat regions and 1.97% consisted of simple sequence repeats (see **Supplementary File S1** for details regarding the repeat types and abundances).

The gene prediction pipeline yielded 63,736 complete genes with start and stop codons and a minimum length of 180 bp. Out of these, 2,472 genes had alternate splice variants. For 86.8% of all genes, a functional annotation could be assigned. Gene density varied widely in the genome, ranging from zero per 100 kb window to 49.7%, with an average and median of 18.2 and 17.6%, respectively. Gene lengths ranged from 180 to 54,183 bp, with an average and median gene length of 3,919 and 3,082 bp, respectively. In the *Fagus sylvatica* genome 4.9 exons per gene were found on average, corresponding well to other high-quality plant genome drafts. The distribution of exons and introns in comparison to *J. regia* and *A. thaliana* is presented in **Table 2**. An analysis of PLAZA genes identified 28,326 such genes in *F. sylvatica*, out of which 1,776 were present in the three other species used for comparison (**Supplementary File S2**).

3.1.2 Telomere and Centromere Predictions

The results given above indicate a high quality of the genome assembly and the gene annotations. To ascertain that the chromosomes were fully resolved, telomeric and centromeric regions were predicted. The tandem repeat element TTAGGG was the most abundant repeat in the genome and was the building block of the telomeric repeats. Out of 12 chromosomes, 8 have stretches of telomeric repeats towards both ends of the chromosomes and, while the other 4 chromosomes have telomeric repeats towards only one end of chromosomes (**Figure 2**). One unplaced scaffold of 110,653 bp which contained 12,051 bp of telomeric repeats at one end, probably represents one of the missing chromosome-ends.

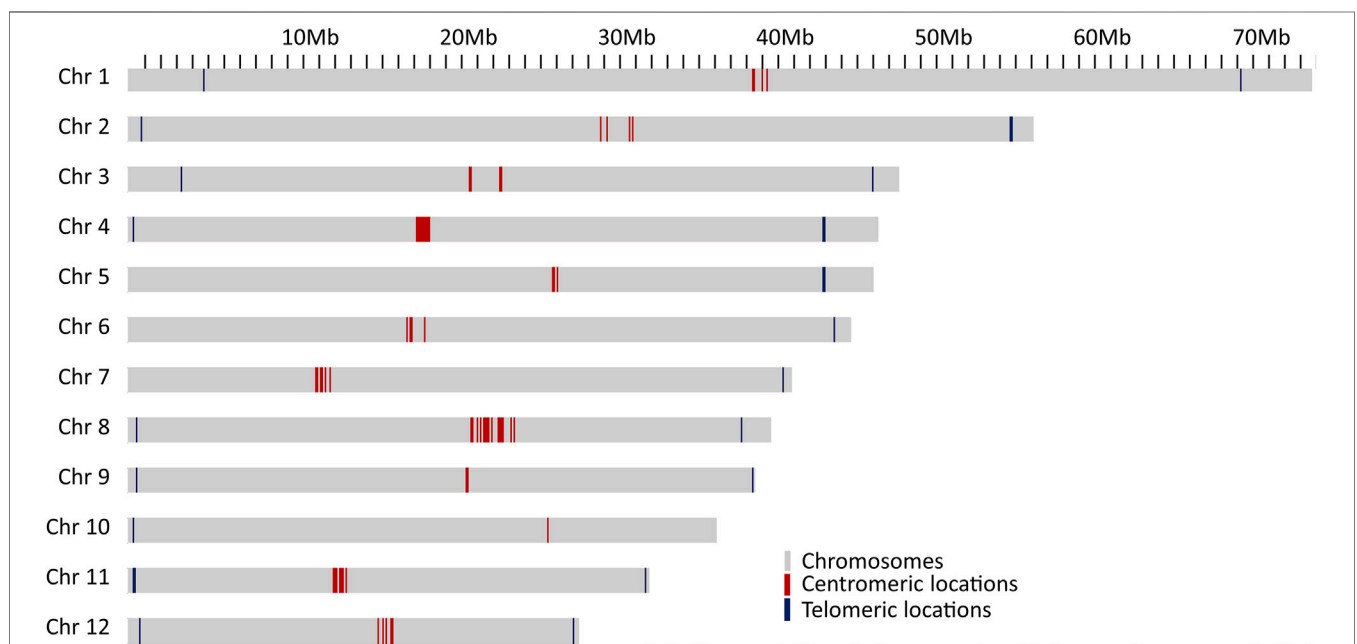
Two different types of potential centromeric repeats were observed, consisting of 79 bp and 80 bp monomer units (**Supplementary File S3**). Centromeric repeats were also observed in higher numbers outside the main centromeric region on several chromosomes (**Supplementary File S3**). However, except for chromosome 10, there was a clear clustering of centromeric repeats within each of the chromosomes, likely corresponding to the actual centromere of the respective chromosomes, and supported also by complementary evidence, such as similarities to centromeric

TABLE 1 | Comparison of BUSCO completeness in Fagaceae genomes available and in the present study (*Fagus sylvatica* V2).

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

TABLE 2 | Distribution of exons in *Fagus sylvatica* in comparison to *Juglans regia* and *Arabidopsis thaliana*.

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

**FIGURE 2** | Locations of probable centromeric repeats on the chromosomes presented as red lines and telomeric locations as blue line on the chromosomes.

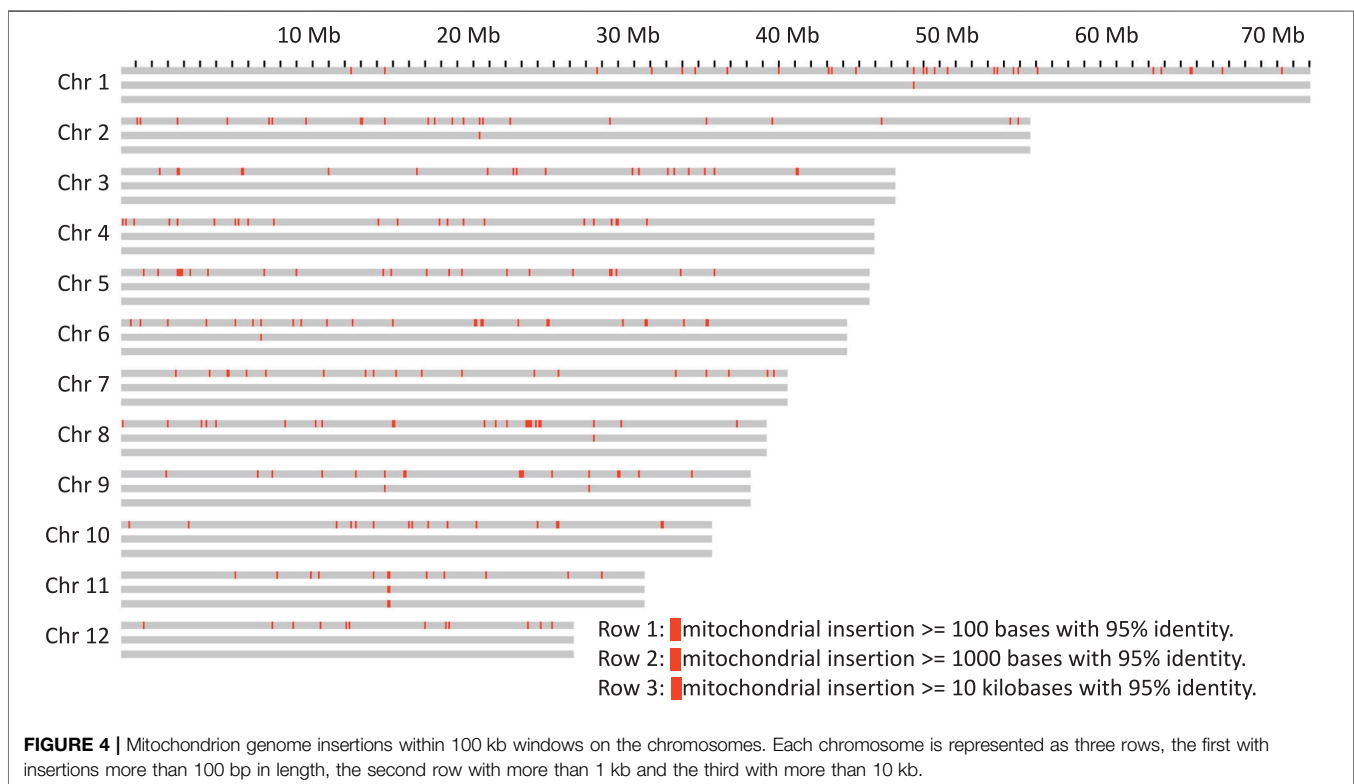
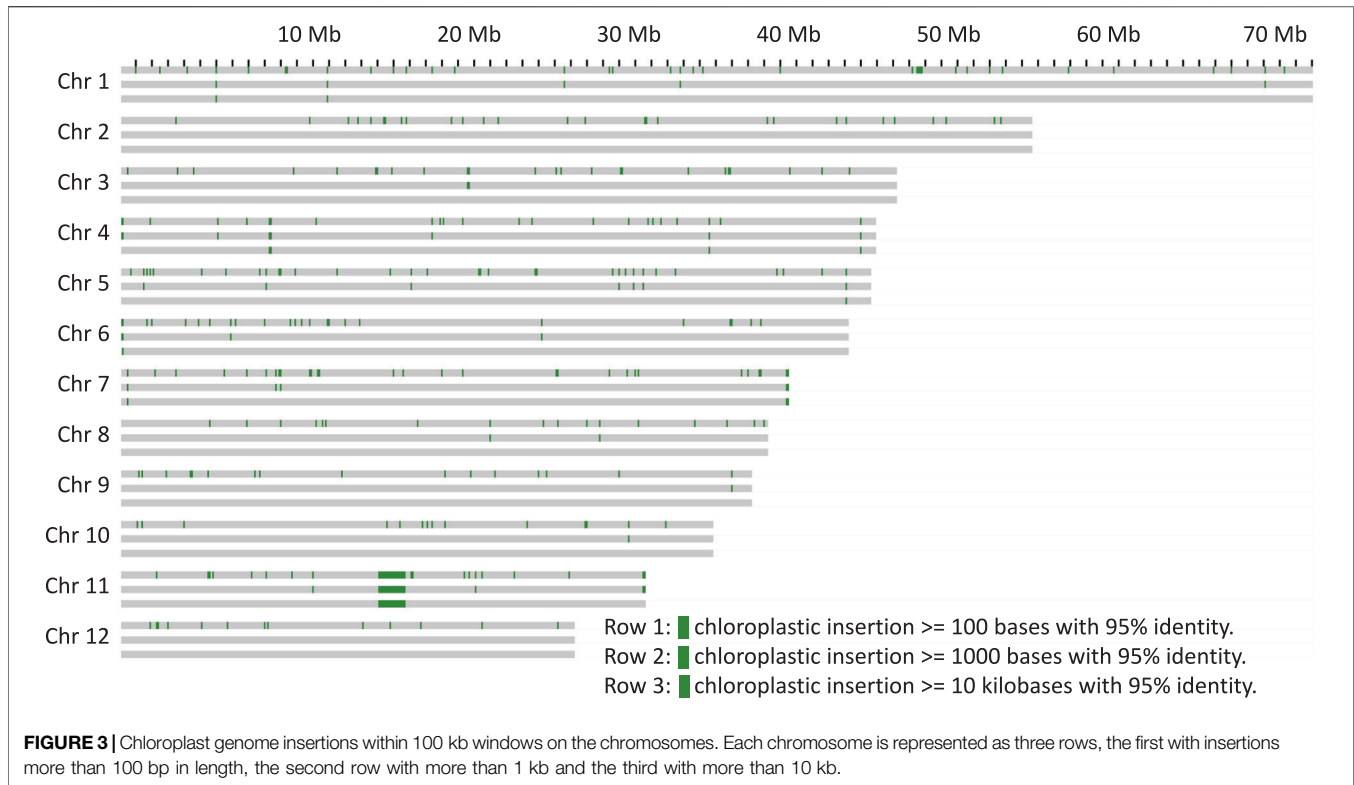
regions of *A. thaliana*, high gypsy element content and low GC content (**Supplementary File S3**).

3.1.3 Integration of Organelle DNA in the Nuclear Genome

As it has previously been shown that organelle DNA insertions can be uneven across the genome and associated with chromatin structure (Wang and Timmis 2013), their distribution in the genome of Bhaga was analysed. For both chloroplast (Mishra et al., 2021a) and mitochondria (Mishra et al., 2021b), multiple integrations of fragments of variable length were observed in all chromosomes (**Figures 3, 4**). These fragments varied in length from the minimum size threshold (100 bp) to 54,784 bp for the

chloroplast and to 26,510 bp for the mitochondrial DNA. The identity of the integrated organelle DNA with the corresponding stretches in the organelle genome ranged from the minimum threshold tested of 95% to 100%. Nuclear-integrated fragments of organelle DNA exceeding 10 kb were found on six chromosomes for the chloroplast, but only on one chromosome for the mitochondrial genome (**Figures 3, 4**).

Nuclear insertions with sequence identity >99% were about ten times more frequent for chloroplast than for mitochondrial DNA with 173 vs 16 for fragments >1 kb and 115 vs 11 for fragments >5 kb, respectively. Eight of these matches of mitochondria were located on unplaced contigs. Overall, mitochondrial insertions tended to be



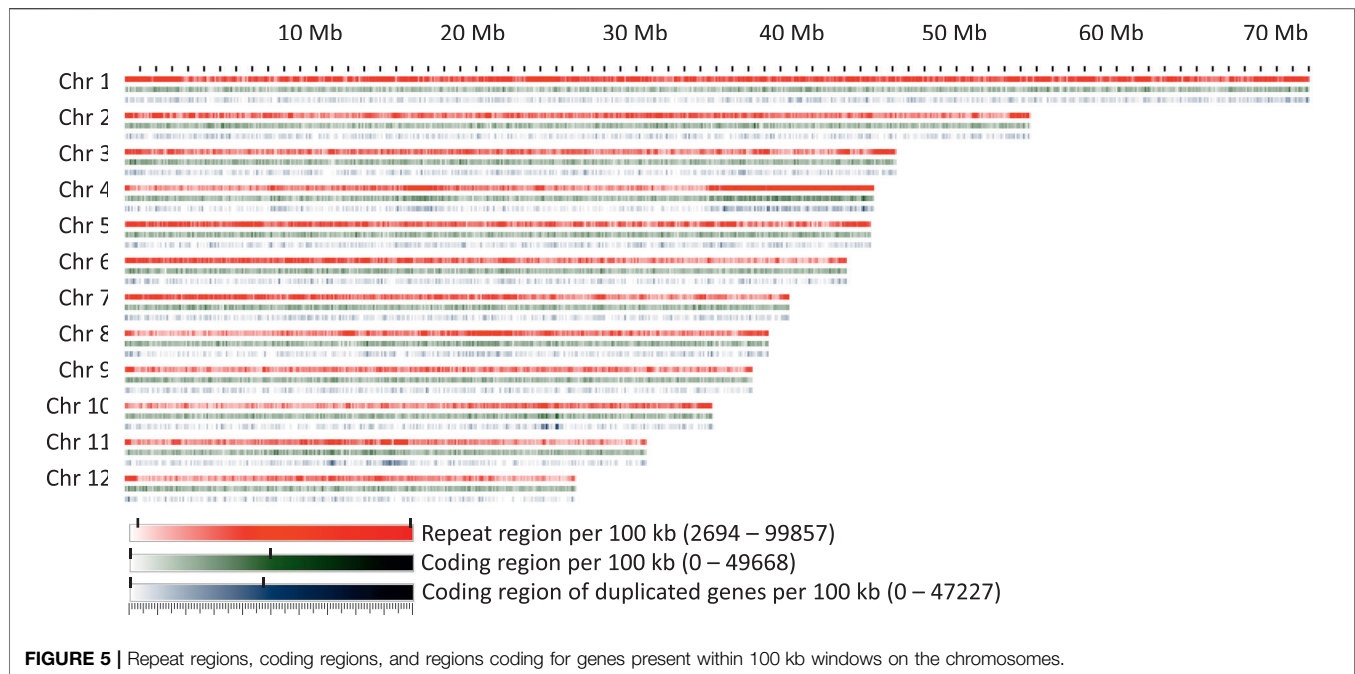


FIGURE 5 | Repeat regions, coding regions, and regions coding for genes present within 100 kb windows on the chromosomes.

smaller and show a slightly higher sequence similarity (**Supplementary File S4**), suggesting that they might be purged from the nuclear genome quicker than the chloroplast genome insertions.

The integration of organelle DNA into the nuclear genome was mostly even, but tandem-like integrations of chloroplast DNA on chromosome 2 were observed (**Figure 3**). In addition, insertions of both organelles were found close to the ends in 4 of the 24 chromosome ends (4, 6, 7, and 8). For the insertions further than 500 kb away from the chromosome ends the integration sites of mitochondrion DNA were sometimes found within the same 100 kb windows where the chloroplast DNA insertion was found. If some regions of the genome are more amenable to the integration of organelle DNA than others needs to be clarified in future studies. A major anomaly was found on chromosome 11, where in a stretch of about 2 Mb (from about Mb 16 to Mb 18 on that chromosome) consisting mainly of multiple insertions of both chloroplast and mitochondrial DNA was observed. In this region, an insertion of more than 20 kb of mitochondrial DNA was flanked by multiple very long integrations of parts of the chloroplast genome on both sides (**Figures 3, 4**). Thus, these integrations appeared almost repeat-like at this particular location.

3.1.4 Repeat Elements and Gene Space

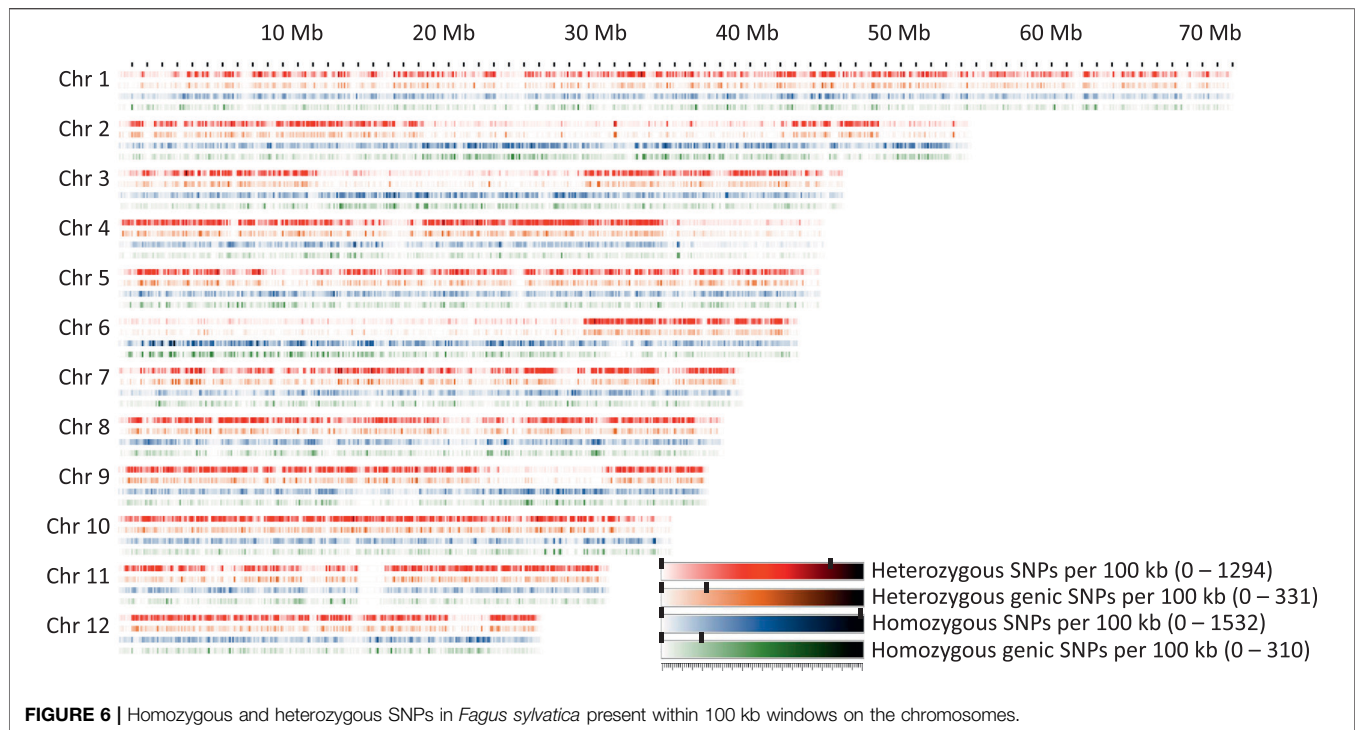
The most abundant repeat elements were LTR elements and LINES, covering 11.49% and 3.66% of the genome, respectively. A detailed list of the element types found, their abundance, and proportional coverage of the genome is given in **Supplementary File S1**. Repeat element presence was variable across the chromosomes (**Figure 5**). While the repeat content per 100 kb window exceeded 50% over more than 88% of

chromosome 1, this was the case for only 37.5% of chromosome 9. Chromosomes showed an accumulation of repeat elements towards their ends, except for chromosome 10, where only a moderate increase was observed on one of the ends, and chromosome 1, where repeat elements were more evenly distributed. Overall, the repeat content was unevenly distributed, with a patchy distribution of repeat-rich and repeat-poor regions of variable length.

A conspicuous anomaly was noticed in chromosome 4, where at one end a large region of about 10 Mb was found in which 97% of the 100 kb windows had a repeat content greater than 70%. This region also contained a high proportion of duplicated or multiplied genes (**Figure 5**). Additional regions containing more than 20% of duplicated genes within a window of at least 1 Mb were identified on chromosomes 4, 10, and 11. On chromosome 11, two clusters were detected, one of which corresponded to the site of organelle DNA insertions described above.

The ribosomal cistrons were reported to be located at the telomeres of four different chromosomes in *F. sylvatica* (Ribeiro et al., 2011). Due to the highly repetitive nature of the ribosomal repeats and their placement near the telomeres, they could not be assigned with certainty to specific chromosomes and thus remained in four unplaced contigs. However, the 5S unit, which is separate from the other ribosomal units in *F. sylvatica*, could be placed near the centromeric locations of chromosomes 1 and 2, in line with the locations inferred by fluorescence microscopy (Ribeiro et al., 2011).

The coding space was more evenly distributed over the chromosomes, with the exception of the regions with high levels of duplicated or multiplied genes. Apart from this, a randomly fluctuating proportion of coding space was observed, with only few regions that seemed to be slightly



enriched or depleted, e.g., in the central part of chromosome 8.

3.1.5 Distribution of Single Nucleotide Polymorphisms

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

Heterozygous SNPs were very unequally distributed over the chromosomes (Figure 6). Several regions, the longest of which comprised more than 30 Mb on chromosome 6, contained only very low amounts of heterozygous SNPs. Apart from the chromosome ends, where generally few heterozygous positions were observed, all chromosomes contained at least one window of 1 Mb where only very few heterozygous SNPs were present. On chromosomes 2, 3, 4, 6, and 9 such areas extended beyond 5 Mb. On chromosome 4 this region corresponded to the repeat region anomaly reported in the previous paragraph, but for the region poor in heterozygous SNPs on chromosome 9, no association with a repeat-rich region could be observed.

Homozygous SNPs differentiating Bhaga and Jamy often followed a different pattern. All regions with low heterozygous SNP frequency longer than 5 Mb had an above-average homozygous SNP frequency, with the exception of the anomalous repeat-rich region on chromosome 4, which had very low frequencies for both homozygous and heterozygous

SNPs. However, there were also two regions of more than 1 Mb length on chromosome 11 that also showed low frequencies of both SNP categories (Figure 6).

Generally, the frequency of overall and intergenic SNPs per 100 kb window corresponded well for both heterozygous and homozygous SNPs, suggesting neutral evolution. However, there were some regions in which genic and intergenic SNP frequencies were uncoupled. For example, on chromosome 1 a high overall heterozygous SNP frequency was observed at 37.7, 48.2, and 56 Mb, but genic heterozygous SNP frequency was low despite normal gene density, suggesting the presence of highly conserved genes. In line with this, also the frequency of homozygous genic SNPs was equally low in the corresponding areas. Similarly, homozygous SNP frequencies were also decoupled on chromosome 1, where a low frequency was observed at 4.2, 7.1, 38.2, 62.1, and 64.8 Mb, but a high genic SNP frequency was observed. This suggests the presence of diversifying genes in the corresponding 100 kb windows, such as genes involved in coping with biotic or abiotic stress.

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

3.1.6 Genome Browser

A genome browser for the genome of Bhaga, with the various genomic features outlined above annotated, is available at

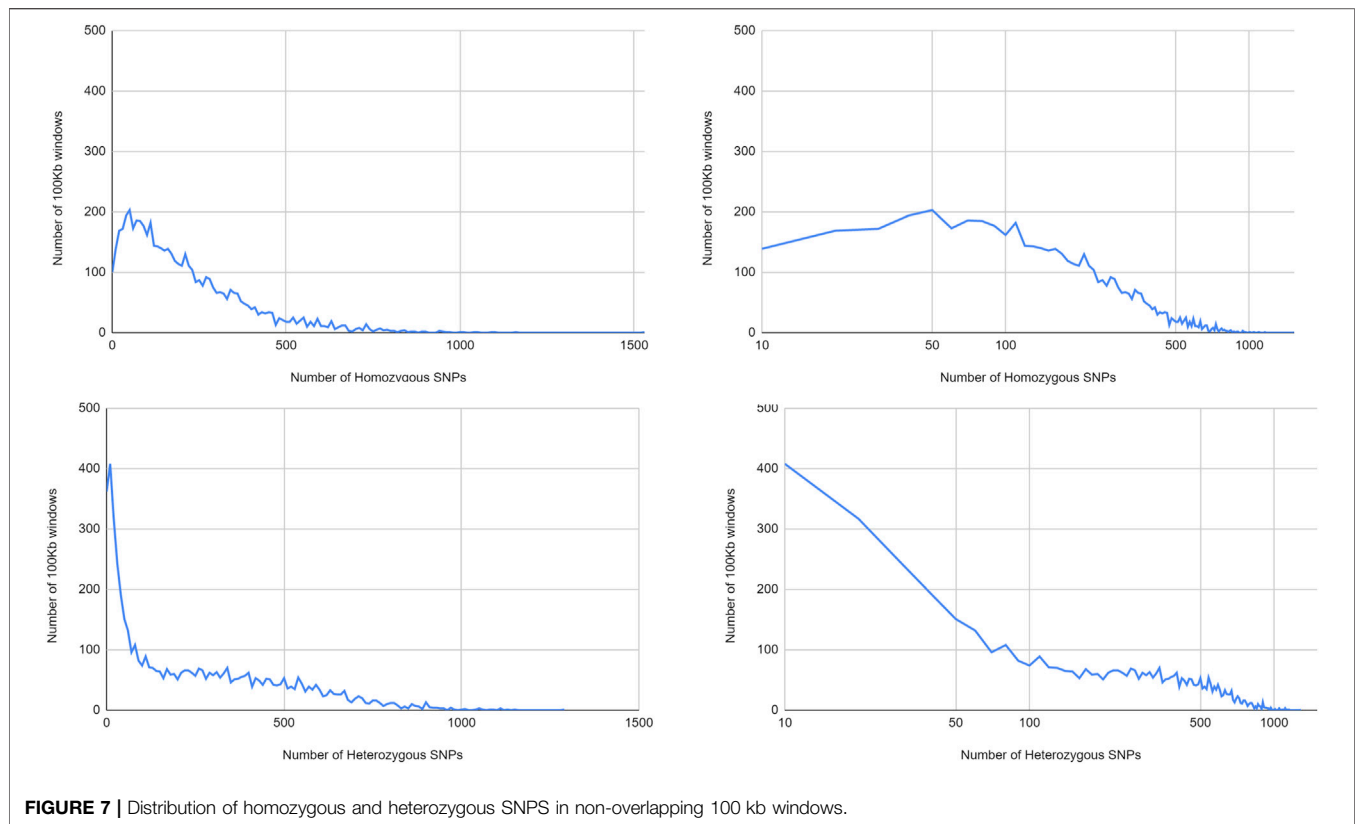


FIGURE 7 | Distribution of homozygous and heterozygous SNPS in non-overlapping 100 kb windows.

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

beechgenome.net. Predicted genes, annotated repeat elements, homozygous SNPs, and heterozygous SNPs are available in “B. Annotations”. The telomeric and centromeric locations, as well as the GC content details are available in “C. Other Details”.

3.2 Validation of Chromosomal-Scale Pseudomolecules

3.2.1 Pedigree Reconstruction

The analysis of the 36 SNPs using Cervus allowed the identification of candidate fathers and reconstruct full-sib families. For 317 of the 537 offspring a likely father was identified. The 19 candidate fathers were represented in the progeny, although their contributions were variable (0.8–21%). For the other offspring, no father could be assigned, i.e., the pollen donor is not present among the surrounding trees (corresponding to 210 genotypes, i.e., 39.1% of the samples when 0 mismatch is allowed, and 22% when 1 mismatch is allowed). The two largest families comprised 68 (MSSBxMSSH) and 86 (MSSBxSSP12) full-sibs. Few years after plantations, 36 genotypes for the former and 49 for the latter survived (Table 3).

3.2.2 A New Unigene Set for European Beech

Our study provides a new reference unigene set for *Fagus sylvatica* based on short and long NGS reads obtained from cDNA libraries constructed from six different tissues. The first unigene set for this species was established back in 2015 using a combination of Sanger and Roche-454 reads (Lesur et al., 2015). The sequences were assembled into 21,000 contigs. A second step was achieved by Müller et al. (2017) using NGS data (Illumina) resulting in 44,000 contigs. This third transcript catalogue contains a total of 34,987 items. When compared to the oak proteome (to date the best annotated among *Fagaceae* species),

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

	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 ^a	Illumina ONT	Velvet Oases	34,987	24,826	22,347
			33,013 ^b (≥200bp)	24,811	21,886

^aIn addition to Illumina and ONT RNAseq, contigs obtained from Lesur et al., 2015 were also included in the analysis. This first unigene provided a total of 609 transcripts to the new reference unigene.

^bTranscripts longer than 200 bp are available online (ENA, accession HBVZ01000000). Smaller contigs are available upon request.

this new reference provides the most complete transcript catalog (Table 4).

3.2.3 Identification of RNAseq-Based SNP Markers for Linkage Mapping

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

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

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

3.2.4 Linkage Map Construction

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

homologous linkage groups obtained from the three different maps (Figure 8).

3.2.5 Alignment of Bhaga Genomic Scaffolds to the SNP-Based Linkage Map of Beech

The 368 mapped markers were aligned on the 12 genomic scaffolds (Bhaga_1 to Bhaga_12) of the *Fagus sylvatica* genome assembly. The alignments were filtered and congruence between scaffolds and linkage groups was checked. Most of the markers from a given LG mapped on a single scaffold (Table 6) providing a genetic validation of the physical assembly obtained for the Bhaga genome sequence. Notable exceptions were: 1) LG11 and LG12, which corresponded to Bhaga_#8; these two chromosomal arms could not be merged into a single LG, and 2) LG13 and scaffold #11, which presented too few markers for unambiguous assignment to one or more scaffolds and LGs, respectively.

4 DISCUSSION

4.1 General Genome Features

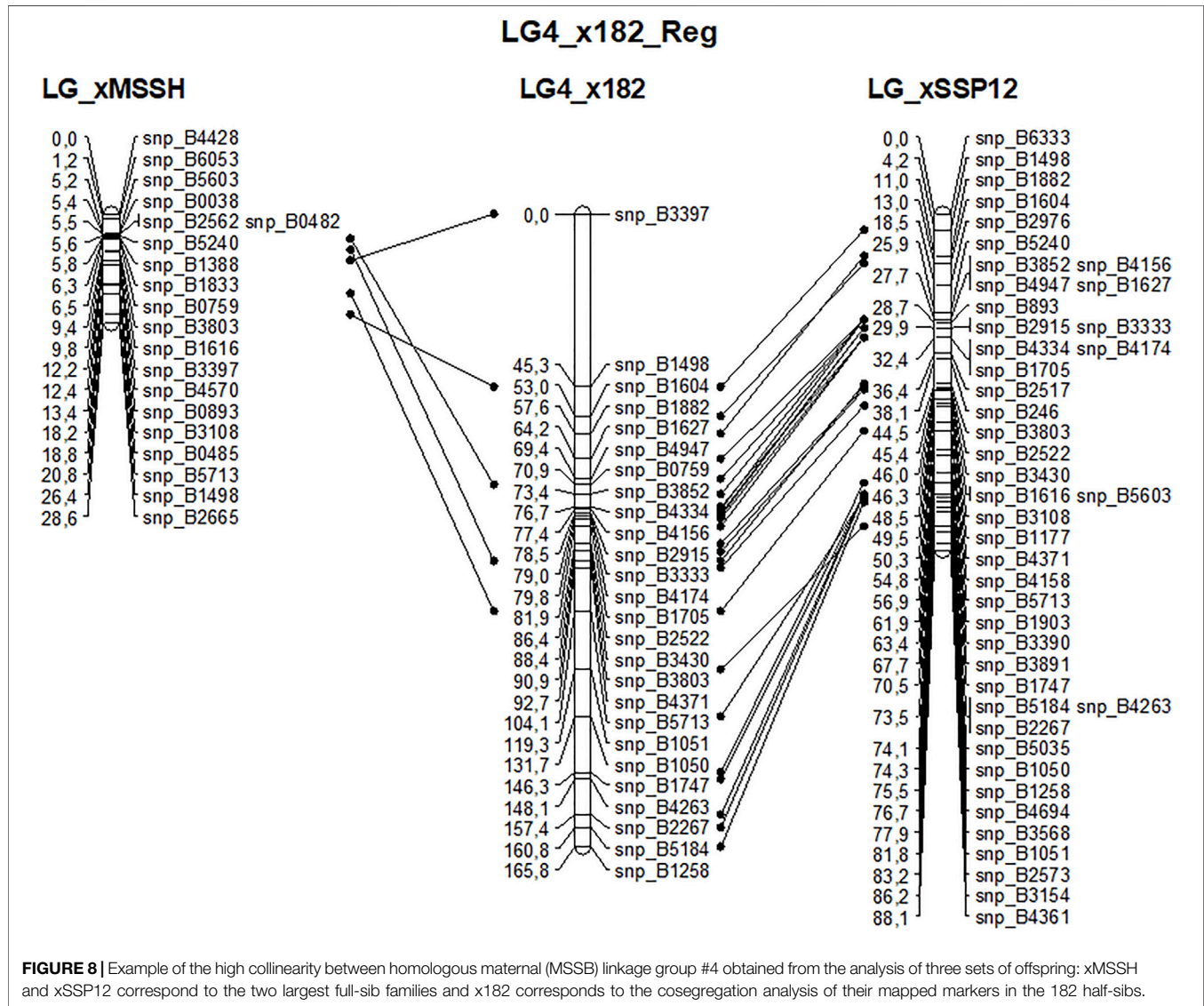
The genome assembled and analysed in this study compares well with previously published *Fagaceae* genomes, both in terms of size and gene space. Here the base chromosome number of 12 could be confirmed, as previously reported based on chromosome counts (Ribeiro et al., 2011). The number of exons per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

4.2 Telomere and Centromere Predictions

Telomeres are inherently difficult to resolve because of long stretches of GC-rich repeats that can cause artefacts during library preparation (Aird et al., 2011) and can lead to biased

TABLE 5 | Characteristics of the combined maternal linkage map in terms of genetic size (cM) and 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	2,213
# of SNPs	37	30	56	36	49	24	24	22	29	15	22	16	8	368



mapping (Dohm et al., 2008). However, using long-read sequencing and Hi-C scaffolding, telomeric repeats could be identified on all chromosomes. It seems likely that several of the unplaced contigs of 4.9 Mb, which included telomeric sequences, were not correctly anchored in the assembly due to ambiguous Hi-C association data resulting from the high sequence similarity of telomeric repeats, because of which for four chromosomes telomeric repeats could be identified only on one of the ends. This might also be due to the presence of ribosomal cistrons on four chromosome ends, which might have interfered with the Hi-C linkage due to their length and

very high sequence similarity. On the outermost regions of the chromosomes, no longer telomeric repeat stretches were present most likely due to their ambiguous placement in the assembly, again because of very high sequence similarity.

Centromere repeats were identified by screening the genome for repeats of intermediate sizes, and were found to be present predominantly within a single location per chromosome. However, lower amounts of centromeric repeat units were also observed to be scattered throughout the genome. The function of the centromeric repeats outside of the centromere remains largely enigmatic but could be associated with chromosome structuring

TABLE 6 | Number of SNP markers of a given linkage group (LG) aligned to a specified scaffold (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	

(Alves et al., 2012) or centromere repositioning (Klein and O'Neill, 2018; Mandáková et al., 2020). Interestingly, two major groups of potential centromeric repeat units of different lengths were found, which did not always coincide. The location of the main occurrence of the centromere-defining repeat unit agreed well with the location previously inferred using chromosome preparations and fluorescence microscopy (Ribeiro et al., 2011).

4.3 Integration of Organelle DNA in the Nuclear Genome

Organelle DNA integration has been frequently found in all kingdoms of life for which high-resolution genomes are available (Stegemann et al., 2003; Guo et al., 2008; Zhang et al., 2020). It can be assumed that this transfer of organelle DNA to the nucleus is the seed of transfer of chloroplast genes to the nuclear genome (Huang et al., 2003). However, apart from a few hints (Yang et al., 2017) it is unclear, which factors stabilise the chloroplast genome so that its content in non-parasitic plants stays relatively stable over long evolutionary timescales (Wang et al., 2007; Xiong et al., 2009). In the present study, it has been found that the insertion of organelle DNA are located mainly in repeat-rich regions of the beech genome. However, their presence in regions without pronounced repeat density might suggest that repeats are not the only factor associated with the insertion of organelle DNA. Nevertheless, it appears that some regions are generally amenable to the integration of organelle DNA, as in several cases chloroplast and mitochondrion insertions were observed in close proximity. The reason for this is unclear, but is known that open chromatin is more likely to accumulate insertions (Wang and Timmis 2013). The potential presence of areas in the genome that are less protected from the insertion of foreign DNA could open up potential molecular biology applications for creating stable transformants.

An anomaly regarding organelle DNA insertion was observed on chromosome 11. Around a central insertion of mitochondrion DNA, multiple insertions of chloroplast DNA were found. The whole region spans more than 2 Mb, which is significantly longer than the organelle integration hotspots reported in other species (Zhang et al., 2020). The evolutionary origin of this large

chromosome region is unclear, but given its repetitive nature it is conceivable that it resulted from a combination of an integration of long fragments and repeat element activity. The presence of multiple copies at the location implies an unusual genome structure in this area, but further analyses, ideally including multiple additional individuals, will be necessary to elucidate the basis for this.

4.4 Distribution of Single Nucleotide Polymorphisms

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

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

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

separated similar stands need to be investigated to disentangle the effects of inbreeding and local adaptation.

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

5 CONCLUSION

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

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

DATA AVAILABILITY STATEMENT

The datasets presented in this study can be found in online repositories. The names of the repository/repositories and accession number(s) are as follows: <https://www.ncbi.nlm.nih.gov/genbank/>, PRJEB24056 and PRJNA450822.

AUTHOR CONTRIBUTIONS

MT conceived the study, with contributions from CP and JB. BU, CB, JB, JM, J-MA, LO, MT, and SP provided materials. All authors conducted laboratory experiments or analysed the data. All authors were involved in data interpretation. BM and MT wrote the manuscript with contributions from the other authors. All authors read and approved the final manuscript.

FUNDING

This study was supported by Grants of the German Science Foundation (Th1632/18-1), National Science Centre, Poland (2012/04/A/NZ9/00F500; 2018/31/F/NZ2/01747), the Polish Ministry of Science, and Higher Education under the program “Regional Initiative of Excellence” in 2019–2022 (Grant No. 008/RID/2018/19), and the LOEWE initiative of the government of Hessen in the framework of the LOEWE Centre for Translational Biodiversity Genomics (TBG). IL received funding from the European Union’s Horizon 2020 research and innovation program (GENTREE project, No. 676876). Further funding was received for the validation of the chromosome-scale assembly, sequencing, computing, and services, respectively: INRAE (Ecodiv division) and the Genoscope: the Commissariat à l’Energie Atomique et aux Energies Alternatives (CEA), 2) France Génomique (ANR-10-INBS-09-08, Beech genome project), Genotoul Bioinformatics Platform Toulouse Occitanie (Bioinfo Genotoul, <https://doi.org/10.15454/1.5572369328961167E12>), 4) the Genome Transcriptome Facility of Bordeaux (grant from Investissements d’Avenir, Convention attributive d’aide EquipEx Xyloforest ANR-10-EQPX-16-01).

ACKNOWLEDGMENTS

We thank many of our colleagues from INRAE (UMR Biogeco) including: Benjamin Dencausse (seed collection in Saint-Symphorien, plantation and tissue sampling in Nouzilly and Guémené-Penfao), Jean-Charles Leplé (bud collection in Nouzilly), Parick Reynet, Edith Reuzeau, Yannick Mellerin (seed collection in Saint-Symphorien), Martine Martin-Clotte (sample preparation for ONF nursery), Maxime Doeland, Jean Broué (DNA extraction and genotyping), Adib Ouayjan (selection of SNPs for pedigree reconstruction), Adline Delcamp, and Emilie Chancerel (mass-array genotyping). We are grateful to the administrative nursery of Guémené-Penfao managed by the Office National des Forêts (ONF) especially Jean-Pierre Huvelin and Olivier Forestier (germination of seeds of the 1st and 2nd campaign, plantation of the 2nd campaign and monitoring of the plantation), as well as the experimental units PAO of INRAE Nouzilly (Benoit Luwez) and GBFOR of INRAE Orléans (Dominique Veisse) for the plantation of the 1st campaign and the monitoring of the plantation.

SUPPLEMENTARY MATERIAL

The Supplementary Material for this article can be found online at: <https://www.frontiersin.org/articles/10.3389/fgene.2021.691058/full#supplementary-material>

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), R18–R14. doi:10.1186/gb-2011-12-2-r18
- Alberti, A., Poulain, J., Engelen, S., Labadie, K., Romac, S., Ferrera, I., et al. (2017). Viral to Metazoan marine Plankton Nucleotide Sequences from the Tara Oceans Expedition. *Sci. Data* 4, 170093. doi:10.1038/sdata.2017.93
- Altschul, S. F., Gish, W., Miller, W., Myers, E. W., and Lipman, D. J. (1990). Basic Local Alignment Search Tool. *J. Mol. Biol.* 215, 403–410. doi:10.1016/s0022-2836(05)80360-2
- Alves, S., Ribeiro, T., Inácio, V., Rocheta, M., and Morais-Cecilio, L. (2012). Genomic Organization and Dynamics of Repetitive DNA Sequences in Representatives of Three *Fagaceae* Genera. *Genome* 55, 348–359. doi:10.1139/g2012-020
- 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. doi:10.1038/nature08800
- Ayad, L. A., and Pissis, S. P. (2017). MARS: Improving Multiple Circular Sequence Alignment Using Refined Sequences. *BMC Genomics* 18 (86), 86–10. doi:10.1186/s12864-016-3477-5
- Bell, M. V. (2018). PLAZA 4.0: an Integrative Resource for Functional, Evolutionary and Comparative Plant Genomics. *Nucleic Acids Res.* (Accessed October 21, 2020).
- Benson, G. (1999). Tandem Repeats Finder: a Program to Analyze DNA Sequences. *Nucleic Acids Res.* 27, 573–580. doi:10.1093/nar/27.2.573
- Bodénès, C., Chancerel, E., Ehrenmann, F., Kremer, A., and Plomion, C. (2016). High-density Linkage Mapping and Distribution of Segregation Distortion Regions in the Oak Genome. *DNA Res.* 23, 115–124. doi:10.1093/dnares/dsw001
- Bolte, A., and Degen, B. (2010). Forest Adaptation to Climate Change - Options and Limitations. *Landbauforsch Volk* 60, 111–117.
- Buchfink, B., Xie, C., and Huson, D. H. (2015). Fast and Sensitive Protein Alignment Using DIAMOND. *Nat. Methods* 12, 59–60. doi:10.1038/nmeth.3176
- Buels, R., Yao, E., Diesh, C. M., Hayes, R. D., Munoz-Torres, M., Helt, G., et al. (2016). JBrowse: a Dynamic Web Platform for Genome Visualization and Analysis. *Genome Biol.* 17 (66), 66–12. doi:10.1186/s13059-016-0924-1
- Ceballos, F. C., Joshi, P. K., Clark, D. W., Ramsay, M., and Wilson, J. F. (2018). Runs of Homozygosity: Windows into Population History and Trait Architecture. *Nat. Rev. Genet.* 19, 220–234. doi:10.1038/nrg.2017.109
- Chen, J. D., Zheng, C., Ma, J. Q., Jiang, C. K., Ercisli, S., Yao, M. Z., et al. (2020). The Chromosome-Scale Genome Reveals the Evolution and Diversification after the Recent Tetraploidization Event in Tea Plant. *Hortic. Res.* 7 (63), 63–11. doi:10.1038/s41438-020-0288-2
- Chen, S., Huang, T., Zhou, Y., Han, Y., Xu, M., and Gu, J. (2017). AfterQC: Automatic Filtering, Trimming, Error Removing and Quality Control for Fastq Data. *BMC Bioinformatics* 18, 80. doi:10.1186/s12859-017-1469-3
- Coumou, D., and Rahmstorf, S. (2012). A Decade of Weather Extremes. *Nat. Clim. Change* 2, 491–496. doi:10.1038/nclimate1452
- Danecek, P., Bonfield, J. K., Liddle, J., Marshall, J., Ohan, V., Pollard, M. O., et al. (2021). Twelve Years of SAMtools and BCFtools. *GigaScience* 10 (2), giab008. doi:10.1093/gigascience/giab008
- Di Filippo, A., Biondi, F., Maugeri, M., Schirone, B., and Piovesan, G. (2012). Bioclimate and Growth History Affect Beech Lifespan in the Italian Alps and Apennines. *Glob. Change Biol.* 18, 960–972. doi:10.1111/j.1365-2486.2011.02617.x
- Dohm, J. C., Lottaz, C., Borodina, T., and Himmelbauer, H. (2008). Substantial Biases in Ultra-short Read Data Sets from High-Throughput DNA Sequencing. *Nucleic Acids Res.* 36, e105. doi:10.1093/nar/gkn425
- Durrant, T. H., De Rigo, D., and Caudullo, G. (2016). “*Fagus sylvatica* in Europe: Distribution, Habitat, Usage and Threats,” in *European Atlas of forest Tree Species*. Editors J. San-Miguel-Ayaz and D. de Rigo (Luxembourg: Publications Office of the European Union), e012b90.
- Flynn, J. M., Hubley, R., Goubert, C., Rosen, J., Clark, A. G., Feschotte, C., et al. (2020). RepeatModeler2 for Automated Genomic Discovery of Transposable Element Families. *PNAS* 117, 9451–9457.
- Flynn, J. M., Lower, S. E., Barbash, D. A., and Clark, A. G. (2018). Rates and Patterns of Mutation in Tandem Repetitive DNA in Six Independent Lineages of *Chlamydomonas reinhardtii*. *Genome Biol. Evol.* 10, 1673–1686. doi:10.1093/gbe/evy123
- Fu, L., Niu, B., Zhu, Z., Wu, S., and Li, W. (2012). CD-HIT: Accelerated for Clustering the Next-Generation Sequencing Data. *Bioinformatics* 28, 3150–3152. doi:10.1093/bioinformatics/bts565
- Geßler, A., Keitel, C., Kreuzwieser, J., Matyssek, R., Seiler, W., and Rennenberg, H. (2007). Potential Risks for European Beech (*Fagus sylvatica* L.) in a Changing Climate. *Trees* 21, 1–11. doi:10.1007/s00468-006-0107-x
- Götz, S., García-Gómez, J. M., Terol, J., Williams, T. D., Nagaraj, S. H., Nueda, M. J., et al. (2008). High-throughput Functional Annotation and Data Mining with the Blast2GO Suite. *Nucleic Acids Res.* 36, 3420–3435.
- Guo, X., Ruan, S., Hu, W., Cai, D., and Fan, L. (2008). Chloroplast DNA Insertions into the Nuclear Genome of Rice: the Genes, Sites and Ages of Insertion Involved. *Funct. Integr. Genomics* 8, 101–108. doi:10.1007/s10142-007-0067-2
- Haas, B. J., Papanicolaou, A., Yassour, M., Grabherr, M., Blood, P. D., Bowden, J., et al. (2013). De Novo transcript Sequence Reconstruction from RNA-Seq Using the Trinity Platform for Reference Generation and Analysis. *Nat. Protoc.* 8, 1494–1512. doi:10.1038/nprot.2013.084
- Ho, E. K. H., Bellis, E. S., Calkins, J., Adrion, J. R., Latta, L. C., and Schaack, S. (2020). Engines of Change: Transposable Element Mutation Rates Are High and Vary Widely Among Genotypes and Populations of *Daphnia magna*. *BioRxiv* [Preprint]. doi:10.1101/2020.09.21.307181
- Hong, Z., Li, J., Liu, X., Lian, J., Zhang, N., Yang, Z., et al. (2020). The Chromosome-Level Draft Genome of *Dalbergia odorifera*. *GigaScience* 9, gaa084. doi:10.1093/gigascience/gaa084
- Huang, C. Y., Ayliffe, M. A., and Timmis, J. N. (2003). Direct Measurement of the Transfer Rate of Chloroplast DNA into the Nucleus. *Nature* 422, 72–76. doi:10.1038/nature01435
- Jiang, S., An, H., Xu, F., and Zhang, X. (2020). Chromosome-level Genome Assembly and Annotation of the Loquat (*Eriobotrya japonica*) Genome. *GigaScience* 9, gaa015. doi:10.1093/gigascience/giaa015
- Jones, P., Binns, D., Chang, H.-Y., Fraser, M., Li, W., McAnulla, C., et al. (2014). InterProScan 5: Genome-Scale Protein Function Classification. *Bioinformatics* 30, 1236–1240. doi:10.1093/bioinformatics/btu031
- Jump, A. S., Hunt, J. M., and Peñuelas, J. (2006). Rapid Climate Change-Related Growth Decline at the Southern Range Edge of *Fagus sylvatica*. *Glob. Change Biol.* 12, 2163–2174. doi:10.1111/j.1365-2486.2006.01250.x
- Kalinowski, S. T., Taper, M. L., and Marshall, T. C. (2007). Revising How the Computer Program CERVUS Accommodates Genotyping Error Increases success in Paternity Assignment. *Mol. Ecol.* 16, 1099–1106. doi:10.1111/j.1365-294x.2007.03089.x
- Kent, W. J. (2002). BLAT—the BLAST-like Alignment Tool. *Genome Res.* 12, 656–664. doi:10.1101/gr.229202
- Kim, D., Langmead, B., and Salzberg, S. L. (2015). HISAT: a Fast Spliced Aligner with Low Memory Requirements. *Nat. Methods* 12, 357–360. doi:10.1038/nmeth.3317
- Klein, S. J., and O’Neill, R. J. (2018). Transposable Elements: Genome Innovation, Chromosome Diversity, and Centromere Conflict. *Chromosome Res.* 26, 5–23. doi:10.1007/s10577-017-9569-5
- Kremer, A., Abbott, A. G., Carlson, J. E., Manos, P. S., Plomion, C., Sisco, P., et al. (2012). Genomics of *Fagaceae*. *Tree Genet. Genomes* 8, 583–610. doi:10.1007/s11295-012-0498-3
- Lalagüe, H., Csilléry, K., Oddou-Muratorio, S., Safrana, J., de Quattro, C., Fady, B., et al. (2014). Nucleotide Diversity and Linkage Disequilibrium at 58 Stress Response and Phenology Candidate Genes in a European Beech (*Fagus sylvatica* L.) Population from Southeastern France. *Tree Genet. Genomes* 10, 15–26. doi:10.1007/s11295-013-0658-0
- Langmead, B., and Salzberg, S. L. (2012). Fast Gapped-Read Alignment with Bowtie 2. *Nat. Methods* 9, 357–359. doi:10.1038/nmeth.1923
- Le Provost, G., Herrera, R., Paiva, J. A., Chaumeil, P., Salin, F., and Plomion, C. (2007). A Micromethod for High Throughput RNA Extraction in Forest Trees. *Biol. Res.* 40, 291–297. doi:10.4067/S0716-97602007000400003

- Lesur, I., Bechade, A., Lalanne, C., Klopp, C., Noirot, C., Leplé, J.-C., et al. (2015). A Unigene Set for European Beech (*Fagus sylvatica* L.) and its Use to Decipher the Molecular Mechanisms Involved in Dormancy Regulation. *Mol. Ecol. Resour.* 15, 1192–1204. doi:10.1111/1755-0998.12373
- Leuschner, C., Meier, I. C., and Hertel, D. (2006). On the Niche Breadth of *Fagus sylvatica*: Soil Nutrient Status in 50 Central European Beech Stands on a Broad Range of Bedrock Types. *Ann. For. Sci.* 63, 355–368. doi:10.1051/forest2006016
- Li, H. (2011). A Statistical Framework for SNP Calling, Mutation Discovery, Association Mapping and Population Genetical Parameter Estimation from Sequencing Data. *Bioinformatics* 27, 2987–2993. doi:10.1093/bioinformatics/btr509
- Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., et al. (2009). The Sequence Alignment/map Format and SAMtools. *Bioinformatics* 25, 2078–2079. doi:10.1093/bioinformatics/btp352
- Li, W., and Godzik, A. (2006). Cd-hit: a Fast Program for Clustering and Comparing Large Sets of Protein or Nucleotide Sequences. *Bioinformatics* 22, 1658–1659. doi:10.1093/bioinformatics/btl158
- Lieberman-Aiden, E., Van Berkum, N. L., Williams, L., Imakaev, M., Ragozcy, T., Telling, A., et al. (2009). Comprehensive Mapping of Long-Range Interactions Reveals Folding Principles of the Human Genome. *Science* 326, 289–293. doi:10.1126/science.1181369
- Ligot, G., Balandier, P., Fayolle, A., Lejeune, P., and Claessens, H. (2013). Height Competition between *Quercus petraea* and *Fagus sylvatica* Natural Regeneration in Mixed and Uneven-Aged Stands. *For. Ecol. Manag.* 304, 391–398. doi:10.1016/j.foreco.2013.05.050
- Madoui, M.-A., Engelen, S., Cruaud, C., Belsler, C., Bertrand, L., Alberti, A., et al. (2015). Genome Assembly Using Nanopore-Guided Long and Error-free DNA Reads. *BMC Genomics* 16, 327. doi:10.1186/s12864-015-1519-z
- Mandáková, T., Hloušková, P., Koch, M. A., and Lysak, M. A. (2020). Genome Evolution in *Arabideae* Was Marked by Frequent Centromere Repositioning. *Plant Cell* 32, 650–665. doi:10.1105/tpc.19.00557
- Marchler-Bauer, A., Lu, S., Anderson, J. B., Chitsaz, F., Derbyshire, M. K., DeWeese-Scott, C., et al. (2011). CDD: a Conserved Domain Database for the Functional Annotation of Proteins. *Nucleic Acids Res.* 39, D225–D229. doi:10.1093/nar/gkq1189
- Marrano, A., Britton, M., Zaini, P. A., Zimin, A. V., Workman, R. E., Puiu, D., et al. (2020). High-quality Chromosome-Scale Assembly of the walnut (*Juglans regia* L.) Reference Genome. *GigaScience* 9, g1aa050. doi:10.1093/gigascience/g1aa050
- Marshall, T. C., Slate, J., Kruuk, L. E. B., and Pemberton, J. M. (1998). Statistical Confidence for Likelihood-based Paternity Inference in Natural Populations. *Mol. Ecol.* 7, 639–655. doi:10.1046/j.1365-294x.1998.00374.x
- Martínez, I., and González-Taboada, F. (2009). Seed Dispersal Patterns in a Temperate forest during a Mast Event: Performance of Alternative Dispersal Kernels. *Oecologia* 159, 389–400. doi:10.1007/s00442-008-1218-4
- Martínez García, P. J., Crepeau, M. W., Puiu, D., Gonzalez-Ibeas, D., Whalen, J., Stevens, K. A., et al. (2016). The walnut (*Juglans regia*) Genome Sequence Reveals Diversity in Genes Coding for the Biosynthesis of Non-structural Polyphenols. *Plant J.* 87, 507–532.
- Michael, T. P., and VanBuren, R. (2020). Building Near-Complete Plant Genomes. *Curr. Opin. Plant Biol.* 54, 26–33. doi:10.1016/j.pbi.2019.12.009
- Mishra, B., Gupta, D. K., Pfenninger, M., Hickler, T., Langer, E., Nam, B., et al. (2018). A Reference Genome of the European Beech (*Fagus sylvatica* L.). *GigaScience* 7, g1y063. doi:10.1093/gigascience/g1y063
- Mishra, B., Ulaszewski, B., Meger, J., Ploch, S., Burczyk, J., and Thines, M. (2021b). A Comparison of Three Circular Mitochondrial Genomes of *Fagus sylvatica* from Germany and Poland Reveals Low Variation and Complete Identity of the Gene Space. *Forests* 12, 571. doi:10.3390/f12050571
- Mishra, B., Ulaszewski, B., Ploch, S., Burczyk, J., and Thines, M. (2021a). A Circular Chloroplast Genome of *Fagus sylvatica* Reveals High Conservation between Two Individuals from Germany and One Individual from Poland and an Alternate Direction of the Small Single-Copy Region. *Forests* 12, 180. doi:10.3390/f12020180
- Mott, R. (1997). EST_GENOME: a Program to Align Spliced DNA Sequences to Unspliced Genomic DNA. *Bioinformatics* 13, 477–478. doi:10.1093/bioinformatics/13.4.477
- Müller, M., Seifert, S., Lübke, T., Leuschner, C., and Finkeldey, R. (2017). De Novo Transcriptome Assembly and Analysis of Differential Gene Expression in Response to Drought in European Beech. *PLoS one* 12 (9), e0184167. doi:10.1371/journal.pone.0184167
- NCBI nr database (2020). <https://ftp.ncbi.nlm.nih.gov/blast/db/> (accessed June 24, 2020).
- Ning, D. L., Wu, T., Xiao, L. J., Ma, T., Fang, W. L., Dong, R. Q., et al. (2020). Chromosomal-level Assembly of *Juglans sigillata* Genome Using Nanopore, BioNano, and Hi-C Analysis. *GigaScience* 9, g1aa006. doi:10.1093/gigascience/g1aa006
- Nong, W., Law, S. T. S., Wong, A. Y. P., Baril, T., Swale, T., Chu, L. M., et al. (2020). Chromosomal-level Reference Genome of the Incense Tree *Aquilaria sinensis*. *Mol. Ecol. Resour.* 20, 971–979. doi:10.1111/1755-0998.13154
- OmicsBox - Bioinformatics Made Easy (2020). BioBam Bioinformatics. <https://www.biobam.com/omicsbox> (Accessed March 3, 2020).
- Ouayjan, A., and Hampe, A. (2018). Extensive Sib-Mating in a Refugial Population of Beech (*Fagus sylvatica*) Growing along a lowland River. *For. Ecol. Manag.* 407, 66–74. doi:10.1016/j.foreco.2017.07.011
- Pfenninger, M., Reuss, F., Kiebler, A., Schönnenbeck, P., Caliendo, C., Gerber, S., et al. (2020). Genomic Basis of Drought Resistance in *Fagus sylvatica*. *BioRxiv*. doi:10.1101/2020.12.04.411264
- Plomion, C., Aury, J.-M., Amsalem, J., Alaïtabar, T., Barbe, V., Belsler, C., et al. (2016). Decoding the Oak Genome: Public Release of Sequence Data, Assembly, Annotation and Publication Strategies. *Mol. Ecol. Resour.* 16, 254–265. doi:10.1111/1755-0998.12425
- Plomion, C., Aury, J.-M., Amsalem, J., Leroy, T., Murat, F., Duplessis, S., et al. (2018). Oak Genome Reveals Facets of Long Lifespan. *Nat. Plants* 4, 440–452. doi:10.1038/s41477-018-0172-3
- Price, A. L., Jones, N. C., and Pevzner, P. A. (2005). De Novo identification of Repeat Families in Large Genomes. *Bioinformatics* 21 Suppl 1 (Suppl. 1_1), i351–8. doi:10.1093/bioinformatics/bti1018
- Priest, S. J., Yadav, V., and Heitman, J. (2020). Advances in Understanding the Evolution of Fungal Genome Architecture. *F1000Res* 9, 776. doi:10.12688/f1000research.25424.1
- Reif, A., Xystrakis, F., Gärtner, S., and Sayer, U. (2017). Floristic Change at the Drought Limit of European Beech (*Fagus sylvatica* L.) to Downy Oak (*Quercus pubescens*) forest in the Temperate Climate of Central Europe. *Not Bot. Horti Agrobot* 45, 646–654. doi:10.15835/nbha45210971
- Rhie, A., Mc Carthy, S. A., Fedrigo, O., Damas, J., Formenti, G., Koren, S., et al. (2020). Towards Complete and Error-free Genome Assemblies of All Vertebrate Species. *BioRxiv* [Preprint]. doi:10.1101/2020.05.22.110833
- Ribeiro, T., Loureiro, J., Santos, C., and Morais-Cecílio, L. (2011). Evolution of rDNA FISH Patterns in the *Fagaceae*. *Tree Genet. Genomes* 7, 1113–1122. doi:10.1007/s11295-011-0399-x
- Rice, P., Longden, I., and Bleasby, A. (2000). EMBOSS: the European Molecular Biology Open Software Suite. *Trends Genet.* 16, 276–277. doi:10.1016/s0168-9525(00)02024-2
- Rose, L., Leuschner, C., Köckemann, B., and Buschmann, H. (2009). Are Marginal Beech (*Fagus sylvatica* L.) Provenances a Source for Drought Tolerant Ecotypes? *Eur. J. For. Res.* 128, 335–343. doi:10.1007/s10342-009-0268-4
- Scaffi, M., Troggo, M., Piovani, P., Leonardi, S., Magnaschi, G., Vendramin, G. G., et al. (2004). A RAPD, AFLP and SSR Linkage Map, and QTL Analysis in European Beech (*Fagus sylvatica* L.). *Theor. Appl. Genet.* 108 (3), 433–441. doi:10.1007/s00122-003-1461-3
- Schulz, M. H., Zerbino, D. R., Vingron, M., and Birney, E. (2012). Oases: Robust De Novo RNA-Seq Assembly across the Dynamic Range of Expression Levels. *Bioinformatics* 28, 1086–1092. doi:10.1093/bioinformatics/bts094
- Seppy, M., Manni, M., and Zdobnov, E. M. (2019). “BUSCO: Assessing Genome Assembly and Annotation Completeness,” in *Methods in Molecular Biology*. Editor M. Kollmar (New York: Humana), Vol. 1962, 227–245. doi:10.1007/978-1-4939-9173-0_14
- Smit, A. F. A., and Hubley, R. (2007). RepeatMasker Open-4.0.5. 2007–2014. <http://www.repeatmasker.org> (Accessed Nov 16, 2020).
- Sork, V. L., Squire, K., Gugger, P. F., Steele, S. E., Levy, E. D., and Eckert, A. J. (2016). Landscape Genomic Analysis of Candidate Genes for Climate Adaptation in a California Endemic oak, *Quercus lobata*. *Am. J. Bot.* 103, 33–46. doi:10.3732/ajb.1500162
- Spinoni, J., Naumann, G., Vogt, J., and Barbosa, P. (2015). European Drought Climatologies and Trends Based on a Multi-Indicator Approach. *Glob. Planet. Change* 127, 50–57. doi:10.1016/j.gloplacha.2015.01.012

- Stanke, M., and Morgenstern, B. (2005). AUGUSTUS: a Web Server for Gene Prediction in Eukaryotes that Allows User-Defined Constraints. *Nucleic Acids Res.* 33 (Suppl. 1_2), W465–W467. doi:10.1093/nar/gki458
- Stegemann, S., Hartmann, S., Ruf, S., and Bock, R. (2003). High-frequency Gene Transfer from the Chloroplast Genome to the Nucleus. *Proc. Natl. Acad. Sci.* 100, 8828–8833. doi:10.1073/pnas.1430924100
- Strijk, J. S., Hinsinger, D. D., Zhang, F., and Cao, K. (2019). *Trochodendron aralioides*, the First Chromosome-Level Draft Genome in *Trochodendrales* and a Valuable Resource for Basal Eudicot Research. *GigaScience* 8, giz136. doi:10.1093/gigascience/giz136
- Tarailo-Graovac, M., and Chen, N. (2009). Using RepeatMasker to Identify Repetitive Elements in Genomic Sequences. *Curr. Prot. Bioinf* 25, 4. doi:10.1002/0471250953.bi0410s25
- Van der Auwera, G. A., and O'Connor, B. D. (2020). *Genomics in the Cloud: Using Docker, GATK, and WDL in Terra*. 1st Edition. Sebastopol, CA: O'Reilly Media.
- Wagner, S., Collet, C., Madsen, P., Nakashizuka, T., Nyland, R. D., and Sagheb-Talebi, K. (2010). Beech Regeneration Research: from Ecological to Silvicultural Aspects. *For. Ecol. Manag.* 259, 2172–2182. doi:10.1016/j.foreco.2010.02.029
- Walker, B. J., Abeel, T., Shea, T., Priest, M., Abouelliel, A., Sakthikumar, S., et al. (2014). Pilon: an Integrated Tool for Comprehensive Microbial Variant Detection and Genome Assembly Improvement. *PLoS ONE* 9, e112963. doi:10.1371/journal.pone.0112963
- Wang, D., and Timmis, J. N. (2013). Cytoplasmic Organelle DNA Preferentially Inserts into Open Chromatin. *Genome Biol. Evol.* 5, 1060–1064. doi:10.1093/gbe/evt070
- Wang, D., Wu, Y.-W., Shih, A. C.-C., Wu, C.-S., Wang, Y.-N., and Chaw, S.-M. (2007). Transfer of Chloroplast Genomic DNA to Mitochondrial Genome Occurred at Least 300 MYA. *Mol. Biol. Evol.* 24, 2040–2048. doi:10.1093/molbev/msm133
- Wang, J. J., Tian, S., Sun, X., Cheng, X., Duan, N., Tao, J., et al. (2020). Construction of Pseudomolecules for the Chinese Chestnut (*Castanea mollissima*) Genome. *G* 10, 3565–3574. doi:10.1534/g3.120.401532
- Wang, L. L., Sun, Y., Sun, X., Yu, L., Xue, L., He, Z., et al. (2020). Repeat-induced point Mutation in *Neurospora crassa* Causes the Highest Known Mutation Rate and Mutational burden of Any Cellular Life. *Genome Biol.* 21 (142), 142–223. doi:10.1186/s13059-020-02060-w
- Xiong, A.-S., Peng, R.-H., Zhuang, J., Gao, F., Zhu, B., Fu, X.-Y., et al. (2009). Gene Duplication, Transfer, and Evolution in the Chloroplast Genome. *Biotechnol. Adv.* 27, 340–347. doi:10.1016/j.biotechadv.2009.01.012
- Yang, F. S., Nie, S., Liu, H., Shi, T. L., Tian, X. C., Zhou, S. S., et al. (2020). Chromosome-level Genome Assembly of a Parent Species of Widely Cultivated Azaleas. *Nat. Commun.* 11 (1), 5269–5313. doi:10.1038/s41467-020-18771-4
- Yang, X., Kang, M., Yang, Y., Xiong, H., Wang, M., Zhang, Z., et al. (2019). A Chromosome-Level Genome Assembly of the Chinese tupelo *Nyssa sinensis*. *Sci. Data* 6 (282), 282–287. doi:10.1038/s41597-019-0296-y
- Yang, X., Yue, Y., Li, H., Ding, W., Chen, G., Shi, T., et al. (2018). The Chromosome-Level Quality Genome Provides Insights into the Evolution of the Biosynthesis Genes for Aroma Compounds of *Osmanthus fragrans*. *Hortic. Res.* 5 (72), 72–13. doi:10.1038/s41438-018-0108-0
- Yang, Z., Hou, Q., Cheng, L., Xu, W., Hong, Y., Li, S., et al. (2017). RNase H1 Cooperates with DNA Gyrase to Restrict R-Loops and Maintain Genome Integrity in *Arabidopsis* Chloroplasts. *Plant Cell* 29, 2478–2497. doi:10.1105/tpc.17.00305
- Ye, C., Hill, C. M., Wu, S., Ruan, J., and Ma, Z. S. (2016). DBG2OLC: Efficient Assembly of Large Genomes Using Long Erroneous Reads of the Third Generation Sequencing Technologies. *Sci. Rep.* 6, 31900–31909. doi:10.1038/srep31900
- Yin, X., Arias-Pérez, A., Kitapci, T. H., and Hedgecock, D. (2020). High-Density Linkage Maps Based on Genotyping-By-Sequencing (GBS) Confirm a Chromosome-Level Genome Assembly and Reveal Variation in Recombination Rate for the Pacific Oyster *Crassostrea gigas*. *G3 - Genes Genom Genet.* 10, 4691–4705. doi:10.1534/g3.120.401728
- Zerbino, D. R., and Birney, E. (2008). Velvet: Algorithms for De Novo Short Read Assembly Using de Bruijn Graphs. *Genome Res.* 18, 821–829. doi:10.1101/gr.074492.107
- Zhang, G.-J., Dong, R., Lan, L.-N., Li, S.-F., Gao, W.-J., and Niu, H.-X. (2020). Nuclear Integrants of Organellar DNA Contribute to Genome Structure and Evolution in Plants. *Int. J. Mol. J. Sci.* 21, 707. doi:10.3390/ijms21030707
- Zhang, X., Zhang, S., Zhao, Q., Ming, R., and Tang, H. (2019). Assembly of Allele-Aware, Chromosomal-Scale Autopolyploid Genomes Based on Hi-C Data. *Nat. Plants* 5, 833–845. doi:10.1038/s41477-019-0487-8
- Zhu, T., Wang, L., You, F. M., Rodriguez, J. C., Deal, K. R., Chen, L., et al. (2019). Sequencing a *Juglans regia* × *J. microcarpa* Hybrid Yields High-Quality Genome Assemblies of Parental Species. *Hortic. Res.* 6 (55), 55–16. doi:10.1038/s41438-019-0139-1

Conflict of Interest: The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

Publisher's Note: All claims expressed in this article are solely those of the authors and do not necessarily represent those of their affiliated organizations, or those of the publisher, the editors and the reviewers. Any product that may be evaluated in this article, or claim that may be made by its manufacturer, is not guaranteed or endorsed by the publisher.

Copyright © 2022 Mishra, Ulaszewski, Meger, Aury, Bodénès, Lesur-Kupin, Pfenninger, Da Silva, Gupta, Guichoux, Heer, Lalanne, Labadie, Opgenoorth, Ploch, Le Provost, Salse, Scotti, Wötzel, Plomion, Burczyk and Thines. This is an open-access article distributed under the terms of the Creative Commons Attribution License (CC BY). The use, distribution or reproduction in other forums is permitted, provided the original author(s) and the copyright owner(s) are credited and that the original publication in this journal is cited, in accordance with accepted academic practice. No use, distribution or reproduction is permitted which does not comply with these terms.