

1 High-quality SNPs from genic regions highlight introgression patterns among
2 European white oaks (*Quercus petraea* and *Q. robur*).

3 *Authors:* Tiange Lang^{1,2,3}, Pierre Abadie^{1,2}, Valérie Léger^{1,2}, Thibaut Decourcelle^{1,2,4}, Jean-
4 Marc Frigerio^{1,2}, Christian Burban^{1,2}, Catherine Bodénès^{1,2}, Erwan Guichoux^{1,2}, Grégoire Le
5 Provost^{1,2}, Cécile Robin^{1,2}, Naoki Tani^{1,2,5}, Patrick Léger^{1,2}, Camille Lepoittevin^{1,2}, Veronica
6 A. El Mujtar^{1,2,6}, François Hubert^{1,2}, Josquin Tibbits⁷, Jorge Paiva^{1,2,8,9}, Alain Franc^{1,2},
7 Frédéric Raspail^{1,2}, Stéphanie Mariette^{1,2}, Marie-Pierre Reviron^{1,2}, Christophe Plomion^{1,2},
8 Antoine Kremer^{1,2}, Marie-Laure Desprez-Loustau^{1,2}, Pauline Garnier-Géré^{1,2,§}

9 Addresses :

10 ¹INRAE, UMR 1202 Biodiversity Genes & Communities, F-33610 Cestas, France

11 ²Univ. Bordeaux, UMR 1202, Biodiversity Genes & Communities, F-33400 Talence, France

12 ³Big Data Decision Institute, Jinan University, Tianhe, Guangzhou, PR China

13 ⁴GEVES, 25 rue Georges Morel, 49071, Beaucozé, France

14 ⁵ Japan International Research Center for Agricultural Sciences (JIRCAS), Tsukuba, Ibaraki,
15 Japan

16 ⁶Unidad de Genética Ecológica y Mejoramiento Forestal. Instituto Nacional de Tecnología
17 Agropecuaria (INTA) EEA Bariloche, Modesta Victoria 4450 (8400), Bariloche, Río Negro,
18 Argentina

19 ⁷ Department of Environment and Primary Industries, Biosciences Research Division,
20 Agribio, 5 Ring Road, Bundoora, Victoria, 3086, Australia

21 ⁸ Instituto de Biologia Experimental e Tecnologica, iBET, Apartado 12, Oeiras 2780-901,
22 Portugal

23 ⁹ Institute of Plant Genetics, Polish Academy of Sciences, 34 Strzeszynska street, Poznan PL-
24 60-479, Poland

25 **Keywords:** SNPs, functional candidate genes, *Quercus robur*, *Q. petraea*, Sanger amplicon
26 resequencing, introgression, species differentiation

27 [§]Corresponding author Pauline Garnier-Géré

28 INRA, UMR 1202 Biodiversity Genes & Communities, F- 33610 Cestas, France; Univ. Bordeaux,
29 UMR 1202 Biodiversity Genes & Communities, Bordeaux, F-33400 Talence, France

30 Fax +33 (0)35385381, email: pauline.garnier-gere@inrae.fr

31 Running title: High-quality SNPs for *Quercus* species

32

33 Abstract

34 In the post-genomics era, non-model species like most *Fagaceae* still lack operational
35 diversity resources for population genomics studies. Sequence data were produced from over
36 800 gene fragments covering ~530 kb across the genic partition of European oaks, in a
37 discovery panel of 25 individuals from western and central Europe (11 *Quercus petraea*, 13
38 *Q. robur*, one *Q. ilex* as an outgroup). Regions targeted represented broad functional
39 categories potentially involved in species ecological preferences, and a random set of genes.
40 Using a high-quality dedicated pipeline, we provide a detailed characterization of these genic
41 regions, which included over 14500 polymorphisms, with ~12500 SNPs -218 being triallelic-,
42 over 1500 insertion-deletions, and ~200 novel di- and tri-nucleotide SSR loci. This catalog
43 also provides various summary statistics within and among species, gene ontology
44 information, and standard formats to assist loci choice for genotyping projects. The
45 distribution of nucleotide diversity ($\theta\pi$) and differentiation (F_{ST}) across genic regions are also
46 described for the first time in those species, with a mean $\theta\pi$ close to ~0.0049 in *Q. petraea*
47 and to ~0.0045 in *Q. robur* across random regions, and a mean F_{ST} ~0.13 across SNPs. The
48 magnitude of diversity across genes is within the range estimated for long-term perennial
49 outcrossers, and can be considered relatively high in the plant kingdom, with an estimate
50 across the genome of 41 to 51 million SNPs expected in both species. Individuals with typical
51 species morphology were more easily assigned to their corresponding genetic cluster for *Q.*
52 *robur* than for *Q. petraea*, revealing higher or more recent introgression in *Q. petraea* and a
53 stronger species integration in *Q. robur* in this particular discovery panel. We also observed
54 robust patterns of a slightly but significantly higher diversity in *Q. petraea*, across a random
55 gene set and in the abiotic stress functional category, and a heterogeneous landscape of both
56 diversity and differentiation. To explain these patterns, we discuss an alternative and non-
57 exclusive hypothesis of stronger selective constraints in *Q. robur*, the most pioneering species
58 in oak forest stand dynamics, additionally to the recognized and documented introgression
59 history in both species despite their strong reproductive barriers. The quality of the data
60 provided here and their representativity in terms of species genomic diversity make them
61 useful for possible applications in medium-scale landscape and molecular ecology projects.
62 Moreover, they can serve as reference resources for validation purposes in larger-scale
63 resequencing projects. This type of project is preferentially recommended in oaks in contrast
64 to SNP array development, given the large nucleotide variation and the low levels of linkage
65 disequilibrium revealed.

66 Introduction

67 High-throughput techniques of the next-generation sequencing (NGS) era and increased
68 genome sequencing efforts in the last decade have greatly improved access to genomic
69 resources in non-model forest tree species (Neale and Kremer 2011, Neale *et al.* 2013;
70 Plomion *et al.* 2016). However, these have only been applied recently to large-scale
71 ecological and population genomics research (Holliday *et al.* 2017). One notable exception
72 are studies undertaken in the model genus *Populus* (e.g. Zhou *et al.* 2014, Geraldès *et al.*
73 2014, Christe *et al.* 2016b) that benefited from the first genome sequence completed in 2006
74 in *P. trichocarpa* (Tuskan *et al.* 2006). In *Fagaceae*, previous comparative mapping and
75 “omics” technologies (reviewed in Kremer *et al.* 2012) with recent development of genomic
76 resources (e.g. Faivre-Rampant *et al.* 2011; Tarkka *et al.* 2013; Lesur *et al.* 2015; Lepoittevin
77 *et al.* 2015, Bodénès *et al.* 2016) set the path to very recent release of genome sequences to
78 the research community (*Quercus lobata*, Sork *et al.* 2016; *Q. robur*, Plomion *et al.* 2016,
79 2018; *Q. suber*, Ramos *et al.* 2018; *Fagus sylvatica*, Mishra *et al.* 2018), and these provide
80 great prospects for future evolutionary genomics studies (Petit *et al.* 2013; Parent *et al.* 2015;
81 Cannon *et al.* 2018; Lesur *et al.* 2018).

82 Recently, building from the European oaks genomic resources (*Quercus Portal* at
83 <https://arachne.pierroton.inra.fr/QuercusPortal/> and references therein), natural populations of
84 4 *Quercus* species (*Q. robur*, *Q. petraea*, *Q. pyrenaica*, *Q. pubescens*) were genotyped for
85 ~4000 single-nucleotide polymorphisms (SNPs, from an initial 8K infinium array, Lepoittevin
86 *et al.* 2015). The data were further analysed (Leroy *et al.* 2017), with results extending
87 previous knowledge on their likely diversification during glacial periods, as well as their
88 recolonization history across Europe and recent secondary contacts (SC) after the last glacial
89 maximum (Hewitt 2000; Petit *et al.* 2002a; Brewer *et al.* 2002). Using recent model-based
90 inference allowing for heterogeneity of migration rates (Roux *et al.* 2014; Tine *et al.* 2014),
91 Leroy *et al.* (2017) showed that the most strongly supported demographic scenarios of species
92 diversification, allowing for gene flow among any pair of the four species mentioned above,
93 included very recent SC, due to a much better fit for patterns of large heterogeneity of
94 differentiation observed across SNP loci (confirmed by Leroy *et al.* 2019, using ~15 times
95 more loci across the genome and the same inference strategy). These recent SC events have
96 been documented in the last decade in many patchily distributed hybrid zones where current
97 *in situ* hybridization can occur among European oak species (e.g. Curtu *et al.* 2007; Jensen *et*
98 *al.* 2009; Lepais and Gerber 2011; Guichoux *et al.* 2013). The resulting low levels of

99 differentiation among *Q. robur* and *Q. petraea* in particular is traditionally linked to a model
100 of contrasted colonization dynamics, where the second-in-succession species (*Q. petraea*) is
101 colonizing populations already occupied by the earlier pioneering *Q. robur* (Petit *et al.* 2003).
102 This model predicts asymmetric introgression towards *Q. petraea* (see Currat *et al.* 2008), as
103 often observed in interspecific gene exchanges (Abbott *et al.* 2003), and a greater diversity in
104 *Q. petraea* was documented at SNP loci showing higher differentiation (Guichoux *et al.*
105 2013). The directionality of introgression in oaks was also shown to depend on species
106 relative abundance during mating periods in particular stands (Lepais *et al.* 2009, 2011).
107 Nevertheless, oaks like other hybridizing taxa are known for the integration of their species
108 parental gene pools and strong reproductive isolation barriers (Muir *et al.* 2000; Muir and
109 Schlötterer 2005; Abadie *et al.* 2012, Lepais *et al.* 2013; Ortiz-Barrientos and Baack 2014;
110 Christe *et al.* 2016a), raising essential questions about the interacting roles of divergent (or
111 other types of) selection, gene flow, and recombination rates variation in natural populations,
112 and their imprints on genomic molecular patterns of variation (e.g. Zhang *et al.* 2016; Christe
113 *et al.* 2016b; Payseur and Rieseberg 2016).

114 These issues will be better addressed with genome-wide sequence data in many samples
115 (Buerkle *et al.* 2011), which will be facilitated in oaks by integrating the newly available
116 genome sequence of *Quercus robur* to chosen HT resequencing methods (Jones and Good
117 2016; e.g. Zhou and Holliday 2012; Lesur *et al.* 2018 for the first target sequence capture
118 study in oaks). However, obtaining high quality haplotype-based data required for nucleotide
119 diversity estimation and more powerful population genetics inferences will likely require the
120 development of complex bioinformatics pipelines dedicated to high heterozygosity genomes
121 and solid validation methods for polymorphism detection (e.g. Geraldès *et al.* 2011; Christe *et*
122 *al.* 2016b).

123 Therefore, the objectives of this work were first to provide a detailed characterization of
124 sequence variation in *Quercus petraea* and *Quercus robur*. To that end, we validated previous
125 unpublished sequence data from the classical Sanger' chain-terminating dideoxynucleotides
126 method (Sanger *et al.* 1977). These sequences targeted fragments of gene regions in a panel of
127 individuals sampled across the western and central European part of both species geographic
128 range. Both functional and expressional candidate genes potentially involved in species
129 ecological preferences, phenology and host-pathogen interactions were targeted, as well as a
130 reference set of fragments randomly chosen across the last oak unigene (Lesur *et al.* 2015).
131 These data were obtained within the framework of the EVOLTREE network activities

132 (<http://www.evoltree.eu/>). Second, we aimed at estimating the distributions of differentiation
133 and nucleotide diversity across these targeted gene regions for the first time in those species,
134 and further test the robustness of comparative diversity patterns observed in the context of
135 both species contrasted dynamics and introgression asymmetry. We discuss the quality,
136 representativity and usefulness of the resources provided for medium scale genotyping
137 landscape ecology projects or as a reference resource for validation purposes in larger-scale
138 resequencing projects.

139 **Material and methods**

140 *Sample collection*

141 The discovery panel (*DiP*) included 25 individuals from 11 widespread forest stands with 2 to
142 4 individuals per location (13 from *Q. robur*, 11 from *Q. petraea*, 1 from *Q. ilex* to serve as
143 outgroup, in Table 1).

144 **Table 1** Geographic location of 25 sampled individuals from *Quercus petraea*, *Q. robur* and
 145 *Q. ilex*.

Country	Sampling site	Latitude	Longitude	Morphological <i>Quercus</i> species	Original Identifier	European cpDNA lineages [#]	cpDNA haplotypes ^{**}		
Spain	Arlaban	42.967	-2.55	<i>petraea</i>	Ar18	B	10, 11, 12		
				<i>robur</i>	Ar22		12		
France	Arcachon	44.663	-1.181	<i>robur</i>	A4*	B	11, 12		
	Pierroton	44.737	-0.776	<i>ilex</i>	IL_C	Euro-Med	H12**		
				<i>robur</i>	11P*			B	10, 12
				<i>robur</i>	3P*				
	Orléans	47.826	1.908	<i>petraea</i>	Qs21*	B	10, 11, 12		
				<i>petraea</i>	Qs28*				
<i>petraea</i>				Qs29*					
Petite Charnie	48.083	-0.167	<i>petraea</i>	PC55	A	7			
			<i>robur</i>	PC229					
			<i>robur</i>	PC233					
Switzerland	Büren	47.105	7.383	<i>petraea</i>	B3	C	1		
				<i>robur</i>	B179				
Hungary	Sopron	47.717	16.642	<i>petraea</i>	S444	A	5, 7		
				<i>robur</i>	S104				
The Netherlands	Meinweg	51.181	6.138	<i>petraea</i>	M51	A, C	1, 5		
				<i>robur</i>	M7				
United Kingdom (UK)	Roudsea Wood	54.218	-3.018	<i>petraea</i>	RW108	B	10, 12		
				<i>robur</i>	RW8				
				<i>robur</i>	RW11				
Germany	Rantzau	53.707	9.765	<i>petraea</i>	R100	A, C	7, 1		
				<i>petraea</i>	R127				
				<i>robur</i>	R300				
				<i>robur</i>	R312				

147 Latitude and longitude are given in the WGS 84 coordinate system. Coordinates correspond either to a
 148 central point in the mixed forest stand, or the mean of individual trees coordinates. *: parents of controlled
 149 crosses used for genetic mapping. [#]: after Petit *et al.* (2002a), the putative glacial refugia for lineage B and
 150 C are located in the south of Spain, and for lineages A and C either in the south of Italy or in the Balkans or
 151 both. **: cpDNA haplotypes are from trees previously sampled in Petit *et al.* (2002b), located within a 50
 152 km radius of studied trees, based on the GD2 database (<http://gd2.pierroton.inra.fr/>). *Quercus* species were
 153 *a priori* assigned from morphological information by persons who sampled the trees, but see below for a
 154 comparison with genetic assignments and introgression analyses of each individual using the STRUCTURE
 155 bayesian inference method (“*Characterization of diversity...*” part).

156 These stands occur across a large part of both *Quercus* species natural distributions, spanning
 157 ~20° in longitude (~2200 km) and ~11° in latitude (~1250 km) in western and central Europe
 158 (Fig. S1, Supporting Information). They are also located in areas covering the three major
 159 cpDNA lineages A, B and C (among five) that indicate different historical glacial refugia

160 (Petit *et al.* 2002a), and extend much further geographically towards northern, eastern and
161 south-eastern European borders (Table 1, after Petit *et al.* 2002b). One stand (Sopron, in
162 Hungary), also occurs within the large geographic distribution of the most Eastern lineage E,
163 in a region where lineages A and C also occur. Individuals were chosen either on the basis of
164 their differing leaf morphology among *Q. robur* and *Q. petraea* species (Kremer *et al.* 2002a),
165 or as parents of mapping pedigrees (e.g. Bodénès *et al.* 2016, see Table 1).

166 Leaves were sampled, stored in silica gel and sent to INRA (Cestas, France) for DNA
167 extraction following Guichoux *et al.* (2013). DNA quality and concentration were assessed
168 with a Nanodrop spectrophotometer (NanoDrop Technologies, Wilmington, 152 DE, USA)
169 and by separating samples in 1% agarose gels stained with ethidium bromide. Extractions
170 were repeated until we obtained at least 20 micrograms of genomic DNA per sample, which
171 was needed for a few thousands individual PCRs.

172 *Choice of genic regions for resequencing*

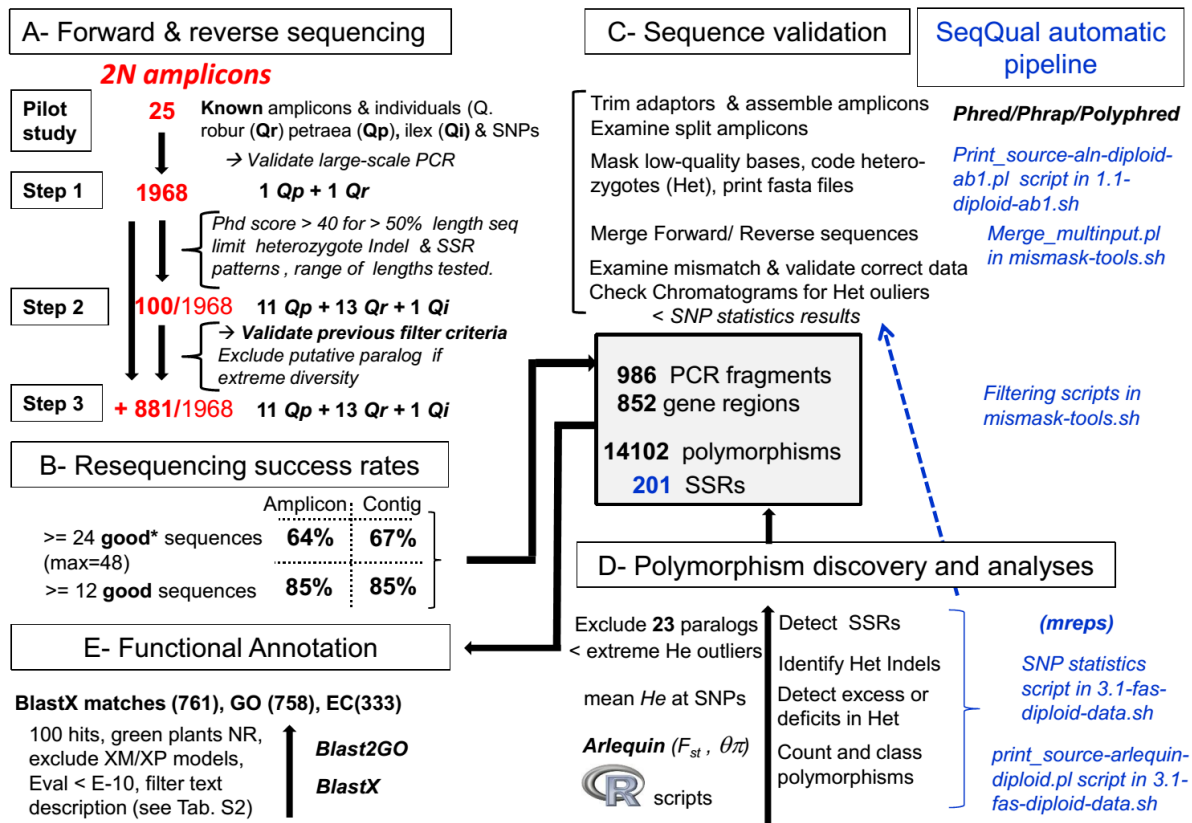
173 Genic regions were chosen from over 103 000 Sanger sequences available in expressed
174 sequence tags (EST) databases at the start of the project. These sequences corresponded to 14
175 cDNA libraries that were prepared with many individuals from both species. They were
176 assembled before finally selecting 2000 fragments for resequencing (Appendix S1 and Fig.
177 S2-A, Supporting information for more on methods producing the original working assembly
178 (*orict*); see also Ueno *et al.* 2010). The targeted fragments were chosen from an extensive
179 compilation of both expressional and functional candidate genes that would likely be involved
180 in white oaks' divergent functions and/or local adaptation, using model and non-model
181 species databases or published results (see Appendix S1 and Fig. S2-B, Supporting
182 information for more details on the strategy followed, and Table S1 for designed primers).

183 *Data production and polymorphism discovery in resequenced fragments*

184 All the sequencing work was performed by Beckman Coulter (Agencourt Bioscience
185 Corporation, Beverly, MA, USA) on ABI3730 capillary sequencers (Applied Biosciences)
186 after preparing DNA samples according to the company's guidelines. Various data quality
187 steps were followed for maximizing the amount and quality of the sequences finally obtained
188 (Fig. 1-A, and Appendix S1, Supporting information for further analyses across 2000
189 amplicons).

190 **Figure 1** Bioinformatics strategy for sequence data production, amplicon assembly,
191 functional annotation, and polymorphism discovery. Scripts used are in italics (see text for

192 further details). GO: Gene Ontology, EC: Enzyme Commission ID. * A **good** sequence is
 193 defined as having a minimum of 50% of its nucleotides with a Phred score above 30.



194

195 Forward and reverse sequences were produced for 986 amplicons across 25 individuals
 196 (100+881 in steps 2 and 3, Fig. 1-A), and more than 85% of them yielded at least 12 high-
 197 quality sequences (Fig. 1-B and column L in Table S1, Supporting information). All amplicon
 198 assembly steps, merging, trimming, and filtering/masking based on quality were performed
 199 with our *SeqQual* pipeline (<https://github.com/garniergere/SeqQual>), with examples of data
 200 and command files. This repository compiles and extends former work dealing with 454 data
 201 (Brousseau *et al.* 2014; El Mujtar *et al.* 2014), providing Bioperl scripts used here that
 202 automatically deal with Sanger haploid or diploid DNA sequences and allow fasta files post-
 203 processing in batch (Fig. 1-C). Sequence variants discovery was finally performed using an
 204 error rate below 0.001 (i.e. Phred score above 30, Ewing *et al.* 1998, and see Appendix S1,
 205 Supporting information for more details). Simple sequence repeat (SSR) patterns were further
 206 detected or confirmed from consensus sequences using the *mreprs* software (Kolpakov *et al.*
 207 2003; see Fig. 1-D, and <https://github.com/garniergere/Reference.Db.SNPs.Quercus/> for a R
 208 script parsing *mreprs* output). Various additional steps involving the treatment of insertion-
 209 deletion polymorphisms (indels) and heterozygote indels (*HI*) in particular, allowed missing

210 data from polymorphic diploid sequence to be minimized (see Appendix S1, Supporting
211 information).

212 *Functional annotation*

213 Resequenced genic regions were annotated using the BlastN best hits of their corresponding
214 *orict* contigs and those of their expected amplicons (*orict-cut*) to most recent oak assembly
215 (*ocv4*, Lesur *et al.* (2015); see Table S2-C, Supporting information). Final consensus
216 sequences for these regions originated from both *orict* and *ocv4* (396 and 368 respectively,
217 see Table S2-A, S2-B, and Appendices S1 and S3, Supporting information), aiming at
218 retrieving the longest sequences, while avoiding to target those with possible chimeric
219 sequences. Functional annotation was then performed via homology transfer using BlastX
220 2.6.0+ program at NCBI (<https://blast.ncbi.nlm.nih.gov/Blast.cgi>) with parameters to optimize
221 speed, hits' annotation description and GO content (Fig. 1-E and Table-S2, Supporting
222 information). Retrieval of GO terms were performed with Blast2GO (Conesa *et al.* 2005 free
223 version at <https://www.blast2go.com/blast2go-pro/b2g-register-basic>) and validation of
224 targeted annotations with Fisher Exact enrichment tests (details of Blast2GO analyses
225 provided in Appendix S1, Supporting information).

226 *Characterization of diversity and genetic clustering*

227 Using the *SNP-stats* script for diploid data from *SeqQual*, simple statistics were computed
228 across different types of polymorphisms (SNPs, indels, SSRs...) including minimum allele
229 frequencies (*maf*) and heterozygote counts, Chi-square tests probability for Hardy-Weinberg
230 proportions, G_{ST} (Nei 1987) and G_{ST}' standardized measure (Hedrick 2005). Complex
231 polymorphisms (involving heterozygote indels (*HI*) and/or SSRs,) were also further
232 characterized (see Appendix S1, Supporting information), and data formatted or analyzed
233 using either Arlequin 3.5 (Excoffier and Lischer 2010), *SeqQual* (e.g. for Arlequin input file
234 with phase unknown, Fig. 1-C), or R scripts. Nucleotide diversity $\theta\pi$ (Nei 1987), based on the
235 average number of pairwise differences between sequences, and its evolutionary variance
236 according to Tajima (1993), were also estimated and compared among species and across
237 candidate genes grouped by broad functional categories (see column F in Table S1,
238 Supporting information), and Weir and Cockerham (1984) F_{ST} estimates of differentiation
239 were computed among species for SNP data along genic regions using analyses of molecular
240 variance (Excoffier 2007).

241 The initial morphological species samples were compared to the genetic clusters obtained
242 with the STRUCTURE v2.3.3 inference method (Falush *et al.* 2003) in order to test possible
243 levels of introgression across individuals. We used the admixture model allowing for mixed
244 ancestry and the correlated allele frequencies assumption for closely related populations as
245 recommended defaults, and since they best represent previous knowledge on each species
246 genetic divergence across their range (e.g. Guichoux *et al.* 2013). Preliminary replicate runs
247 using the same sample of loci produced very low standard deviation across replicates of the
248 data log likelihood given K ($\ln \Pr(X/K)$, see Fig. S3-A, Supporting information). We thus
249 resampled loci at random for each of 10 replicate datasets in 3 different manners to add
250 genetic stochasticity: 1) one per region, 2) one per 100 bp block, and 3) one per 200 bp block
251 along genes (see Appendix S1, Supporting information and
252 <https://github.com/garniergere/Reference.Db.SNPs.Quercus/tree/master/STRUCTURE.files>
253 for examples of STRUCTURE files as recommended by Gilbert *et al.* (2012), along with R
254 scripts for outputs). Statistical independence among loci within each species was verified with
255 Fisher's exact tests implemented in Genepop 4.4 (Rousset 2008).

256 **Results**

257 *Polymorphisms typology and counts*

258 Among the amplicons tested, 986 were successful, 13 did not produce any data and 23 were
259 excluded because of paralog amplifications (Fig. 1-C and Table S1, Supporting information).
260 Around 25% of the successful amplicons overlapped and were merged, consistently with their
261 original design across contigs. Despite the presence of *HI* patterns due to SSR or indels, most
262 amplicons were entirely recovered with forward and reverse sequencing. Several (5% of the
263 total) were however kept separate, either because of functional annotation inconsistency, or
264 because amplicon overlap was prevented by the presence of SSRs or putative large introns
265 (see "Final gene region ID" column with -F/-R suffix in Table S1, Supporting information).
266 We finally obtained 852 genic regions covering in total ~529 kilobases (kb), with an average
267 size of 621 bp per region, ranging from 81 to 2009 bp (Table 2, and Appendix S4, Supporting
268 information, for genomic consensus sequences).

269 **Table 2** Typology of polymorphisms in successfully resequenced amplicons.

	Both species and introgressed individuals	<i>Q. petraea</i>	<i>Q. robur</i>	<i>Q. ilex</i>
Total length resequenced (bp)	529281	-	-	196676
Number (Nb) of amplicons	986	-	-	486
Nb of genic regions	852	-	-	394
Mean genic region size - N50 size (bp)	621-700	-	-	500-539
Minimum - Maximum genic region size (bp)	81-2009	-	-	198-1285
Estimated intron sequences (bp)	186827	-	-	-
Mean haploid sample size (total sequence)	34.71	13.35	18.28	-
Polymorphism in 852 genic regions				
Mean haploid sample size (variants)	32.16	12.57	13.85	-
Monomorphic genic regions	15 (1.76%)	18 (2.14%)	21 (2.52%)	-
Genes with at least one single base indel	591	345	379	-
" " " " one larger indel (>1 bp)	252	190	214	-
" " " " one SSR (>=di)	163	-	-	-
SNPs only (excluding 1 bp indels)	12478	7511	8078	-
Indels (1 bp)	1213	751	809	-
Indels (2-5 bp)	221	142	161	-
Indels (6-10 bp)	88	72	71	-
Indels (11-50 bp, excl. SSRs)	98	81	79	-
Indels (74,146,219,341 bp, excl. SSRs)	4	3	4	-
Total number of polymorphisms	14102	8560	9202	676
<i>Triallelic SNPs</i>	<i>218</i>	<i>141</i>	<i>165</i>	-
<i>...Singletons (incl. 1 bp indels)</i>	<i>4334</i>	<i>1990</i>	<i>2151</i>	-
<i>...Variable SSRs (excl. homopolymers)</i>	<i>111</i>	-	-	-
Total length with sequence variant positions	17594	10765	11451	-
Sequence length of indels and complex polymorphisms (Indels and SSRs)	5116	-	-	-

270 Counts for *Q. petraea* exclude the 2 most introgressed individuals (Qs28 and S444 in Table 1); SSR: simple
 271 sequence repeats; "N50 size" is the size for which the cumulative sum of gene amplicons' size equal or
 272 higher than this value corresponds to 50% of the total amplicons' size sum; The number of polymorphisms
 273 for *Q. ilex* equals the number of heterozygotes in the resequenced individual across amplicons; Numbers of
 274 monomorphic regions were computed for those with at least 10 gametes in both species; Some detected
 275 SSR patterns were not polymorphic in our samples (detailed in Tables S1 and S5, Supporting information).

276 Compared to the EST-based expected total fragment size of ~ 357 kb, around 187 kb of intron
 277 sequence was recovered across 460 of the resequenced regions (assuming intron presence if
 278 an amplicon size was above its expected size by 40 bp). Introns represented ~35% of genic
 279 regions in length and ~51% of those including introns.

280 We observed 14102 polymorphisms in both species across 852 gene regions, 15 of those
 281 regions (<2%) being monomorphic (Table 2). This corresponds to 1 polymorphism per ~38
 282 bp, or 1 per ~30 bp when considering the total number of variant positions in both species
 283 (17594 bp, Table 2). Remarkably, variant positions involving larger indels, SSRs and mixed

284 complex polymorphism patterns represented ~30% of the total variant positions (Table 2, and
285 see their exhaustive lists with various statistics in Table S3 and S4, Supporting information).
286 We observed 12478 SNPs (88.5% of all polymorphisms), 1 SNP per 42 bp, and 218 triallelic
287 SNPs (~1.75% of SNPs) were confirmed by visual examination of chromatograms.

288 Considering only one species, we observed on average 1 variant position per ~48 bp, 1
289 polymorphism per ~60 bp, and 1 SNP per ~68 bp. Among indels, 1213 (8.6% of all
290 polymorphisms) were single base, 309 ranged from 2 to 10 bp, and 102 had sizes above 10 bp
291 which were mostly shared among species (Table 2). In this range-wide sample, there were
292 4334 singletons among all single base polymorphisms, 506 of them being indels. Overall,
293 indels were present in 69% of gene regions and non-single base ones across ~30% of them.
294 Excluding homopolymers (see Appendix S1, Supporting information), we detected 201 SSRs
295 occurring on 163 gene regions by considering a minimum repeat numbers of 4 and a
296 mismatch rate among repeats below 10% (Table 2, Table S1 and Table S5, Supporting
297 information), and 55% (111) were polymorphic in our sample of individuals (Table 2).
298 Among them, 89 (44%) had dinucleotide repeats and 65 (32%) trinucleotide repeats. The
299 SSRs with the lowest number of repeats (<5) had a majority (59%) of repeat sizes between 4
300 and 7, the rest being trinucleotides (Table S5, Supporting information).

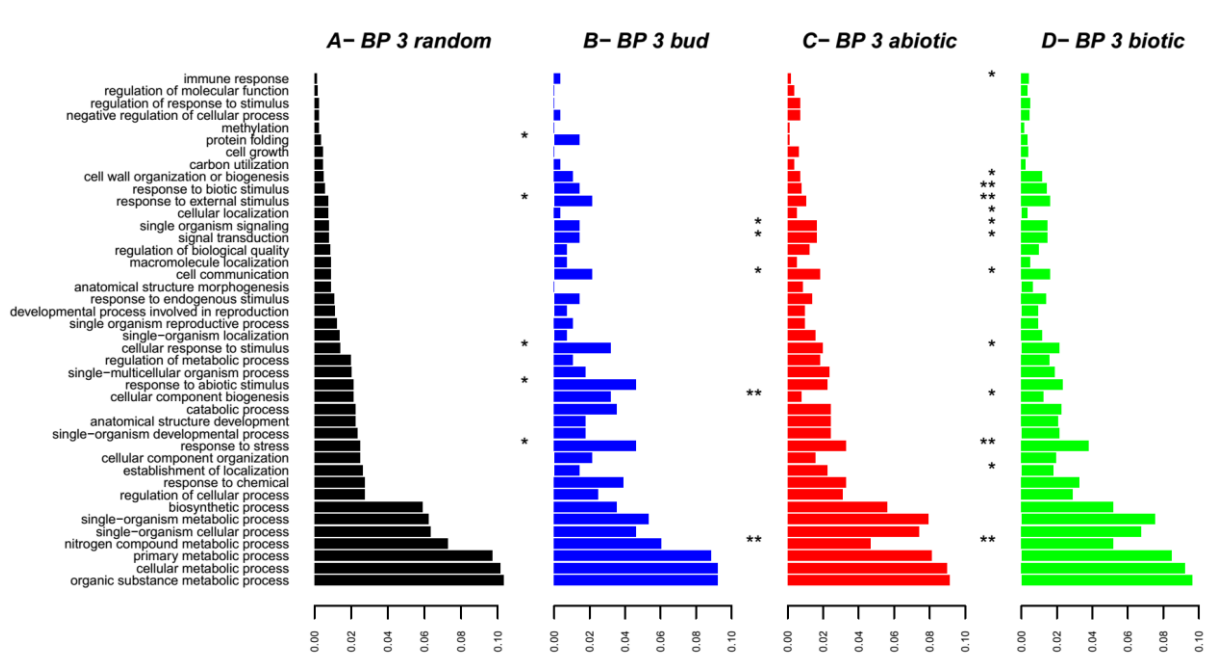
301 Using the same PCR conditions, homologous sequence data were obtained for one individual
302 of the outgroup *Quercus ilex* across 37% of the gene regions (~197 kb, 397 sequences, 676
303 heterozygous sites in Table 2), which illustrates both their sequence similarity yet divergence
304 for a species belonging to the *Ilex* versus *Quercus* taxonomic group (Lepoittevin *et al.* 2015;
305 see Table S1 column Q, and see Appendix S5, Supporting information, for *Q. ilex* genomic
306 sequences).

307 *Annotations and GO term distributions*

308 BlastX matches with *E*-values below 10^{-30} were found for ~97% (738/764) of the contig
309 consensus, only 11 sequences (1.4%) having hits with *E*-values above 10^{-10} that were all
310 among the reference random sample (see BlastX criteria in Table S2, Supporting
311 information). The most represented species among the best hits with informative annotations
312 were *Prunus persica* (111), *Theobroma cacao* (91), *Morus notabilis* (57) and *Populus*
313 *trichocarpa* (45) (Appendix S6-A, Supporting information), which probably illustrates both
314 the close phylogenetic relationships among *Quercus* and *Prunus* genera, consistently with
315 results obtained on the larger *ocv4* assembly (Lesur *et al.* 2015), and the quality and
316 availability of *P. persica* genome annotation (Verde *et al.* 2013, 2017).

317 Between 1 to 30 GO terms could be assigned to 761 sequences, with EC codes and
 318 InterProScan identifiers for 343 and 733 of them respectively (Fig. 1, and Table S2,
 319 Supporting information). The most relevant GO terms were then retained using the Blast2GO
 320 “annotation rule” (Conesa *et al.* 2005) that applies filters from the Direct Acyclic Graph
 321 (DAG) at different levels (Fig. 2, Fig. S4-A- to-F, Supporting information).

322 **Figure 2** Distributions of GO terms across different gene lists (*bud*, *abiotic* and *biotic*) at
 323 biological process level 3, and Fisher exact tests across pairs of sequence clusters with the
 324 same GO terms between the random list and other lists. Significance levels *: P<0.05, **: P<0.01.
 325



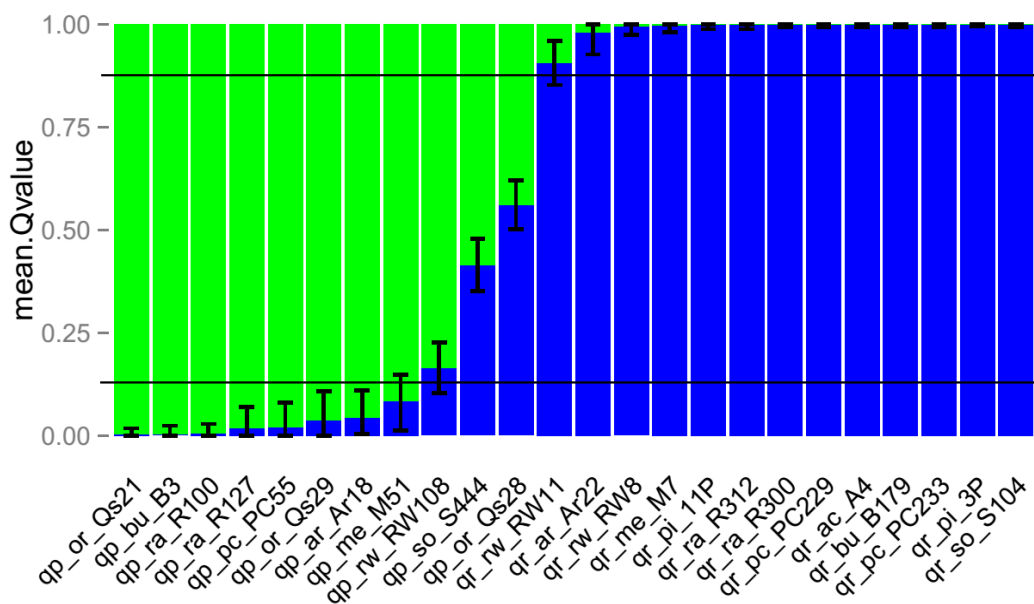
326
 327 At biological process (BP) level 3, apart from general terms involving “metabolic processes”,
 328 a large number of sequences (between ~100 and ~150) were mapped to “response to...” either
 329 “...stress”, “...abiotic stimulus” or “...chemical”, and also to categories linked to
 330 developmental processes (Fig. S4-D, Supporting information).

331 Enrichment tests also revealed a significant increase at both BP levels 2 and 3 for the
 332 following GO categories: “response to stress” or “external stimulus” for *bud* and *biotic* gene
 333 lists, “response to abiotic stimulus” for the *bud* list, and “immune” and “biotic stimulus”
 334 responses for the *biotic* list (see Fig. 2-B to 2-D compared to Fig. 2-A, and Fig. S5,
 335 Supporting information). Most of these exact tests (>80%) were still significant when
 336 selecting genes attributed exclusively to one particular list (in Table S1, Supporting
 337 information), which adds to the relevance of our original gene lists in targeting particular
 338 functional categories.

339 *Species assignment and introgressed individuals*

340 In both species, the proportion of significant association tests among the loci used for
 341 clustering (> two million within each species) was generally one order of magnitude below
 342 the type-I error rates at 5% or 1%. This indicates a very low background LD within species at
 343 their range levels, consistently with the underlying model assumptions used in STRUCTURE.
 344 Based on both $\ln \text{Pr}(X/K)$ and ΔK statistics and as expected, the optimal number of genetic
 345 clusters inferred was 2, whatever the number of polymorphisms and type of sampling (Fig. 3,
 346 Fig. S3 and S6, Supporting information).

347 **Figure 3** Posterior assignment probabilities of individuals into two optimal clusters from
 348 STRUCTURE analyses, sorted in increasing order of belonging to cluster 2 (here *Q. robur* (Qr,
 349 in blue/dark grey), the alternative cluster 1 matching *Q. petraea* (Qp, in green/light grey),
 350 apart from individuals with higher introgression levels. Each bar represents one individual and
 351 includes mean upper and lower bounds of 90% Bayesian confidence intervals around mean *Q*-
 352 values across 10 replicates. Each replicate is a different random sample of 1785
 353 polymorphisms. Horizontal black lines represent the 0.125 and 0.875 values, which can be
 354 considered as typical thresholds for back-crosses and later-generation hybrids (Guichoux *et*
 355 *al.* 2013), values within those thresholds suggesting a mixed ancestry with the other species
 356 for a small number of generations in the past.



357
 358 Most individuals (20) clearly belonged to either cluster with a mean probability of cluster
 359 assignment above 0.9, which was not significantly different from 1, based on mean values of
 360 90% Bayesian credible intervals (BCI) bounds across replicates, and for different types of
 361 sampling or SNP numbers (Fig. 3 and Fig. S6, Supporting information). Two individuals from
 362 Roudsea Wood in UK, the most northerly forest stand of this study, were considered to be
 363 significantly introgressed, each from a different cluster, since both showed a BCI that did not

364 include the value “1” across other replicated runs and SNP sampling (Fig. S6, Supporting
365 information), RW108 also having a mean probability above 0.125 (Fig. 3). Although M51 has
366 a mean assignment value close to that of RW11 in the particular run shown in Fig.3, its BCI
367 was larger and often included the zero value in other runs (Fig. S6, Supporting information),
368 so it was assigned to the *Q. petraea* cluster. In the initial morphological *Q. petraea* group, two
369 individuals were also clearly of recent mixed ancestry: one from the easternmost forest stand
370 of Sopron (S444), and another one (Qs28) from central France, considered previously to be a
371 *Q. petraea* parental genotype in two oak mapping pedigrees (Bodénès *et al.* 2012, 2016;
372 Lepoittevin *et al.* 2015). However, Qs28 shows here a clear F1 hybrid pattern, given its
373 probability values close to 0.5 and its BCI maximum upper and minimum lower bound values
374 of 0.30 and 0.61 respectively across runs (Fig. 3 and Fig. S6-A to S6-J, Supporting
375 information). Testing 3 or 4 possible clusters showed the same ancestry patterns for the
376 introgressed individuals with 2 main clusters and similar *Q*-values (data not shown), which
377 does not support alternative hypotheses of introgression from different species in those
378 individuals.

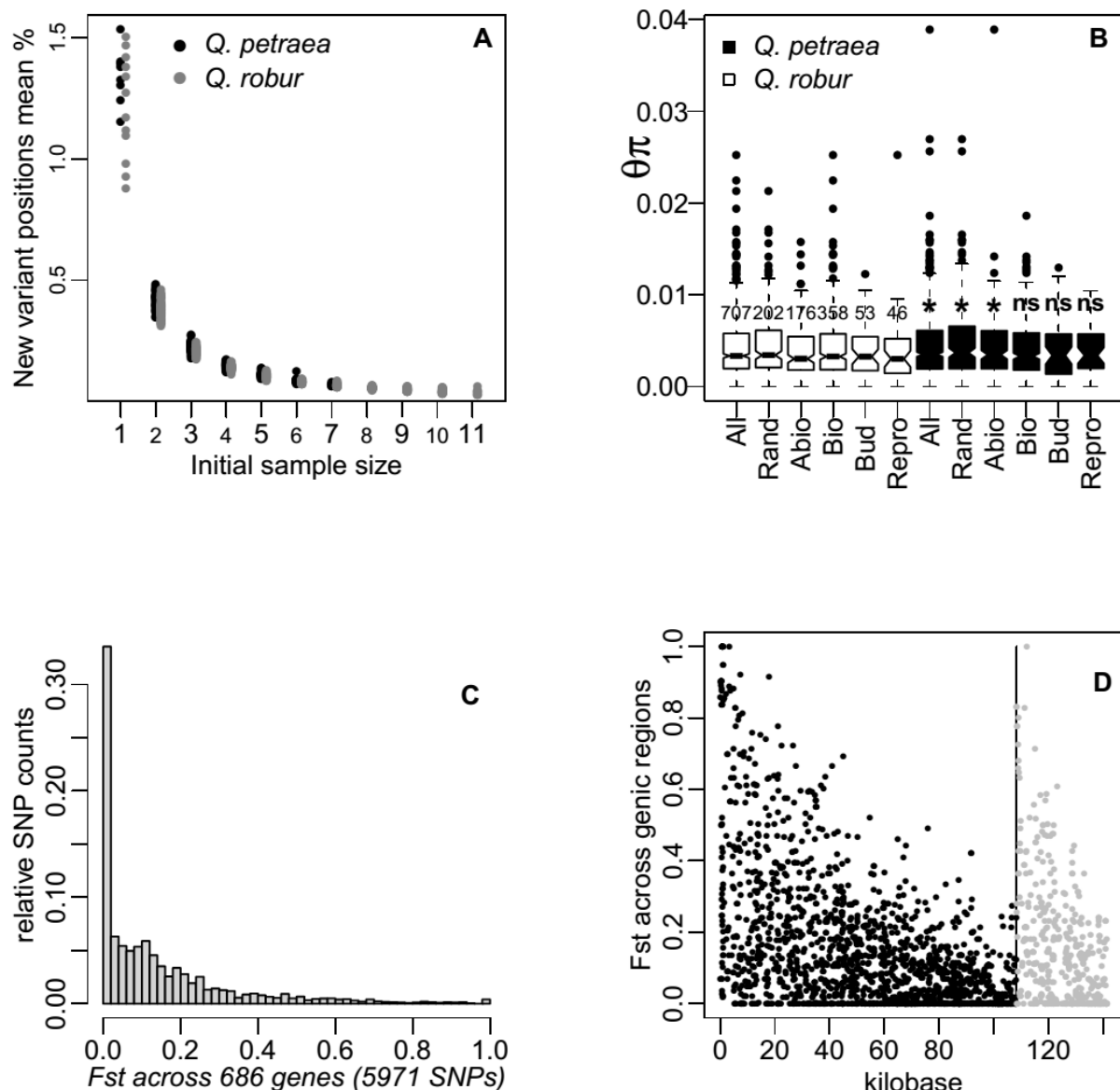
379 *Large heterogeneity of diversity and differentiation across genes*

380 Nucleotide diversity was thus estimated in each parental species after excluding Qs28,
381 RW108, S444 and RW11, which were considered to be the 4 most introgressed individuals
382 (see Fig. 3 above). We then checked how the remaining samples represented species'
383 diversity. Starting with one individual, we observe a dramatic drop in the mean proportion of
384 new variant positions brought by each new individual in any species (*Mpn*) as a function of
385 the initial sample size, followed by a subsequent stabilization (Fig. 4-A, and see Fig. S7-A,
386 Supporting information). Indeed, *Mpn* was only around 11% when going from 4 to 5
387 individuals in both species, and stabilized below 5% after 8 individuals in *Q. robur* (Fig. 4-A).
388 We thus decided to retain 726 gene regions with at least 8 gametes per species (listed in
389 column L in Table S1, Supporting information). The larger *Q. robur* sample after excluding
390 the most introgressed individuals (24 versus 16 gametes in *Q. petraea*) only exhibited
391 slightly higher polymorphism counts than in *Q. petraea* overall (Table 3).

392 Also, 48% and 52% of the polymorphisms observed were exclusive to *Q. petraea* and *Q.*
393 *robur* respectively in our panel, the rest being shared among species (Table 3). Among
394 exclusive polymorphisms, 46% and 44% were singletons in *Q. petraea* and *Q. robur*
395 respectively, suggesting that they might be either rare in both species, or more polymorphic in
396 local populations from which few individuals were sampled across the species wider ranges.

397 Overall and within both species, we observed a large variation in the numbers of segregating
 398 sites per gene size (Fig. S7-B, Supporting information).

399 **Figure 4** Mean proportion of new variant sites brought by each new distinct individual added
 400 to all possible initial sample size combinations (A); Mean nucleotide diversity (considering all
 401 polymorphisms) in both species across genic regions, and different functional categories (B)
 402 compared between species with Wilcoxon signed-rank tests: significant at $P < 5\%$ (*), non-
 403 significant (ns); Histogram of F_{st} estimates across polymorphic gene regions with a minimum
 404 of 8 gametes per species, after excluding singletons and grouping negative with null values
 405 (C); Manhattan plot of F_{st} estimates sorted by mean F_{st} values across randomly chosen
 406 (black dots) and Bud phenology (grey dots) genic regions (D).



407
 408 The mean nucleotide diversity estimates ($\theta\pi$) across genic regions when considering all
 409 polymorphisms were 0.00447 and 0.00425 in *Q. petraea* and *Q. robur* respectively, with up to

410 a 10-fold variation among polymorphic genes overall and in different functional categories
 411 (Fig. 4-B and Table 3).

412 **Table 3** Polymorphism counts and nucleotide diversity in parental species across genic
 413 regions with larger sample sizes.

Polymorphism in 726 gene fragments	both species	<i>Q. petraea</i>	<i>Q. robur</i>
Number of individuals considered	20	8	12
Monomorphic gene fragments	17 (2.34%)	19 (2.63%)	20 (2.87%)
Total number of polymorphisms	11089	7061	7721
SNPs only	9867	6226	6830
All Indels and SSRs	1222	835	891
Exclusive polymorphisms	-	3359	4024
Singletons among them (%)	-	0.456	0.437
Shared polymorphisms	3696	-	-
Mean nucleotide diversity estimates*			
<i>SNPs only</i>	3.849E-03	3.957E-03**	3.740E-03
" " diversity range		0-0.03823	0-0.02525
Tajima's evolutionary standard deviation	2.549E-03	2.632E-03	2.465E-03
SNPs only (509 chosen genes)	3.752E-03	3.821E-03	3.682E-03
SNPs only (202 random genes)	4.103E-03	4.306E-03	3.900E-03
<i>All polymorphisms</i>	4.359E-03	4.471E-03	4.247E-03
" " diversity range		0-0.03893	0-0.02525
Tajima's evolutionary standard deviation	2.816E-03	2.903E-03	2.729E-03
All polymorphisms (509 chosen genes)	4.214E-03	4.278E-03	4.150E-03
All polymorphisms (202 random genes)	4.716E-03	4.944E-03	4.488E-03

414 The 4 most introgressed individuals from Fig. 3 (Qs28, S444, RW108, RW11) are excluded for
 415 computations. Monomorphic regions are defined as in Table 2. *: Diversity is computed for regions
 416 with a minimum of 200 bp overall and at least 8 gametes per species at variant positions. The 509
 417 chosen genes belong to the different functional categories listed in Table S1. Values in the "both
 418 species" column for diversity estimates are means across all genes, of both species' values. **: Values
 419 in bold indicate significant Wilcoxon paired ranked tests for a higher *Q. petraea* nucleotide diversity
 420 compared to *Q. robur* across genes.

421 When including SNPs only, mean $\theta\pi$ decreased overall by more than 10% (Table 3, and see
 422 column D in Table S4, Supporting information). The large variation among genes is also
 423 illustrated by the absence of significant differences between mean diversity among functional
 424 categories *within species*, in most comparisons using non-parametric Wilcoxon rank sum tests
 425 (*Wrs*) with similar number of genes. Two notable exceptions were observed when considering
 426 all polymorphisms: the *biotic stress* category (358 genes) had on average a lower $\theta\pi$ in *Q.*
 427 *petraea* than in the random gene list (211 genes, *Wrs* Pr<0.042), and the mean $\theta\pi$ of the
 428 *reproductive phenology* category was significantly lower in both species than that of the *Bud*
 429 *phenology* category (*Wrs* Pr<0.040 and Pr<0.013 in *Q. petraea* and *Q. robur* respectively,

430 considering exclusive categories from Table S2, Supporting information). Genes with $\theta\pi$
431 estimates above 0.02 were found across most categories, whether considering all
432 polymorphisms (Fig. 4-B) or SNPs only. The 8 genic regions showing the highest $\theta\pi$ values
433 in both species were annotated for example as disease resistance, transcription factor or
434 membrane transport proteins, half of them being from the original random list.

435 Comparing nucleotide diversity between individuals according to their main cpDNA lineages
436 B versus A or C (Table 1), no significant differences were found between lineages within both
437 species, using *Wpr* tests across all genes (see also the lineage-associated distributions of
438 genes' diversity in Fig. S8, Supporting information). This was also true for all functional
439 categories. In both species, the mean differentiation across genes among lineages was very
440 low (<0.015 , each gene estimate being the mean F_{ST} across all polymorphisms at this gene),
441 with very few genes (~1%) having much higher mean F_{ST} (ranging from 0.21 to 0.41 or 0.56
442 within *Q. petraea* and *Q. robur* respectively).

443 Mean $\theta\pi$ comparison tests *between species* across all gene regions were not significant (Table
444 3, *Wrs* $Pr > 0.15$ for all polymorphisms or SNPs only), nor were they across different
445 categories and between gene pairs, using a 95% confidence interval based on Tajima's
446 evolutionary variance for $\theta\pi$ (Tajima 1983) while assuming underlying Gaussian
447 distributions. Indeed for the same genic regions, many examples can be found of higher $\theta\pi$
448 estimates in one species or the other. However, comparing diversity estimates across the exact
449 same positions and performing Wilcoxon paired ranked tests (*Wpr*) across all genes, there was
450 a significant pattern of a slightly higher diversity in *Q. petraea* (see Table 3 and Fig. 4-B),
451 whether considering all polymorphisms (*Wpr* $Pr < 0.028$) or SNPs only (*Wpr* $Pr < 0.036$). This
452 pattern remained significant across the 202 polymorphic genes chosen randomly (*Wpr*
453 $Pr < 0.037$, all polymorphisms, Table 3), even when excluding the 5% or 10% of genes having
454 the highest $\theta\pi$ values. This pattern of a significantly higher $\theta\pi$ in *Q. petraea* was not
455 observed when considering the 509 polymorphic gene regions chosen in functional categories,
456 either together or separately in the different categories (Fig. 4-B), except for the *Abiotic stress*
457 category.

458 We also observed a very large variation for F_{ST} estimates across gene regions and functional
459 categories, which covered the full range of possible values [0,1], with mean values of ~0.13
460 whether considering all polymorphisms or SNPs only (Fig. 4-C, and Fig. 4-D for the random
461 genic regions and a representative example in one category). The very few segregating sites

462 with F_{ST} values equal to one had either missing individuals' or strands, possibly caused by
463 polymorphisms within primer regions. Among the sites sequenced for the full sample of
464 gametes, the 20 highest F_{ST} values ranged from 0.6 to 0.9 and belonged to 10 genic regions,
465 many of which also showed null or very low F_{ST} values within 100 bp. This large variation in
466 differentiation was observed between very close variant sites in many genes, suggesting very
467 high recombination rates at genome-wide and range-wide scales, and consistently with the
468 very low expected background LD (see above). Additionally, a large variance is expected
469 around F_{ST} estimates due to the relatively low sample size in both species, in particular for bi-
470 allelic loci (Weir and Hill 2002; Buerkle *et al.* 2011; e.g. Eveno *et al.* 2008).

471 **Discussion**

472 In the NGS era, non-model tree species such as many *Fagaceae* still lag behind model species
473 for easy access to sequence polymorphism and SNP data (but see Gugger *et al.* 2016 for
474 *Quercus lobata*). These data are needed for larger scale studies addressing the many diversity
475 issues raised by their combined economic, ecological and conservation interests (Cavender-
476 Bares 2016; Fetter *et al.* 2017; Holliday *et al.* 2017). Recent achievements and data
477 availability from the *Q. robur* genome sequence project (Plomion *et al.* 2018) opens a large
478 range of applications in many related temperate and tropical *Fagaceae* species due to their
479 conserved synteny (Cannon *et al.* 2018). In this context, we discuss below the representativity
480 of our data in terms of species genomic diversity as well as the robust patterns observed
481 across genes, and further illustrate their past and future usefulness for *Quercus* species.

482 *Genic resources content, quality, and representativity*

483 We provide a high-quality polymorphism catalog based on Sanger resequencing data for more
484 than 850 gene regions covering ~530 kb, using a discovery panel (*DiP*) from mixed *Q. robur*
485 and *Q. petraea* populations located in the western and central European part of their
486 geographic range. This catalog details functional annotations, previous published information,
487 allele types, frequencies and various summary statistics within and across species, which can
488 assist in choosing novel polymorphic sites (SNPs, SSRs, indels...) for genotyping studies.
489 Among genomic SSRs, more than 90% (~200) are new (17 already detected in Durand *et al.*
490 2010; 3 in Guichoux *et al.* 2011), so they constitute an easy source of potentially polymorphic
491 markers in these oak species. Standard formats for high-density genotyping arrays and primer
492 information are also provided, making these resources readily operational for medium scale
493 molecular ecology studies while avoiding the burden of bioinformatics work needed for SNP
494 development (Tables S1 to S5, Supporting information, and see also

495 <https://github.com/garniergere/Reference.Db.SNPs.Quercus> for additional information). This
496 catalog corrects and largely extends the SNP database for *Q. petraea/robur* at
497 <https://arachne.pierroton.inra.fr/QuercusPortal/> which was previously used to document a
498 SNP diversity surrogate for both *Quercus* species in the oak genome first public release
499 (Plomion *et al.* 2016).

500 Thanks to a high quality dedicated pipeline, we could perform a quasi-exhaustive
501 characterization of polymorphism types in our *DiP* and across part of the genic partition of
502 these *Quercus* species (see Fig. 1). Although base call error rates below 1/1000 were used (as
503 originally developed for Sanger sequencing), most variant sites were located in regions with
504 lower error rates (below 1/10000) so that true singletons could be identified. At the genotypic
505 level, a Sanger genotyping error rate below 1% was previously estimated using a preliminary
506 subset of around 1200 SNPs from this catalog (corresponding to around 5800 data points in
507 Lepoittevin *et al.* 2015). This rate can be considered as an upper bound for the present study,
508 given all additional validation and error correction steps performed. Although little produced
509 now with the advent of NGS methods, Sanger data have served for genome sequencing
510 projects in tree species before 2010 (Neale *et al.* 2017), and have been instrumental, in
511 combination to NGS for BAC clones sequencing, in ensuring assembly long-distance
512 contiguity in large genomes such as oaks (Faivre-Rampant *et al.* 2011, Plomion *et al.* 2016).
513 Sanger sequencing has also provided reference high-quality data to estimate false discovery or
514 error rates, and validate putative SNPs in larger scale projects (e.g. Geraldès *et al.* 2011 in
515 *Populus trichocarpa*; Sonah *et al.* 2013 in Soybean; Cao *et al.* 2014 in *Prunus persica*).

516 Finding an optimal balance between the number of samples and that of loci is critical when
517 aiming to provide accurate estimates of diversity or differentiation in population genetics
518 studies. Given the increasing availability of markers in non-model species (usually SNPs), it
519 has been shown by simulation (Willing *et al.* 2012, Hivert *et al.* 2018) and empirical data
520 (Nazareno *et al.* 2017) that sample sizes as small as 4 to 6 individuals can be sufficient to
521 infer differentiation when a large number of bi-allelic loci (> 1000) are being used. A broad-
522 scale geographic sampling is however required if the aim is to better infer genetic structure
523 and complex demographic scenarios involving recolonization and range shifts due to past
524 glacial cycles, such as those assumed for many European species (Lascoux and Petit 2010,
525 Keller *et al.* 2010, Jeffries *et al.* 2016, Sousa *et al.* 2014). Our sampling design is likely to
526 have targeted a large part of both species overall diversity and differentiation across the
527 resequenced genic regions. This is first suggested by the small proportion of additional

528 polymorphisms once an initial sample of 8 gametes was included for each species (i.e. ~10%
529 and decreasing as sample size increases, Fig. 4-A and Fig. S7-A, Supporting information).
530 Considering the *DiP* within each species, each individual brings on average ~166 new
531 variants (~1% of the total). Second, the large variance observed across gene nucleotide
532 diversity estimates (see Table 3) is mostly due to stochastic evolutionary factors rather than to
533 sampling effects so unlikely to be impacted by sample sizes over 10 gametes (Tajima 1983).
534 Third, sampling sites are located in regions which include 4 out of the 5 main cpDNA
535 lineages reflecting white oaks recolonization routes (lineages A to C and E in Petit *et al.*
536 2002a), the likely haplotypes carried by the *DiP* individuals being A to C (Table 1).

537 Therefore, if new populations were being sampled within the geographical range considered,
538 they would likely include many of the alleles observed here within species and at other genes
539 across their genomes. For differentiation patterns, older and more recent reports showed a low
540 genetic structure among distant populations within each species, and a relatively stable overall
541 differentiation among species compared to possible variation across geographical regions
542 (Bodénès *et al.* 1997; Mariette *et al.* 2002; Petit *et al.* 2003; Muir and Schlotterer 2005;
543 Derory *et al.* 2010; Guichoux *et al.* 2013; Gerber *et al.* 2014). For new populations sampled
544 outside the *DiP* geographic range, a recent application to *Q. robur* provenances located in the
545 low-latitude range margins of the distribution (where 3 main cpDNA lineages occur) showed
546 a high rate of genotyping success, a high SNP diversity, and outliers potentially involved in
547 abiotic stress response (Temunovic *et al.* 2020).

548 We further tested the frequency spectrum representativity of our range-wide *DiP* by
549 comparing genotypic data for a set of 530 independent SNPs (called *sanSNP* for Sanger data)
550 with data for the same set of SNPs obtained in Lepoittevin *et al.* (2015, called the *illuSNP* set
551 since it used the Illumina Infinium array technology) for larger numbers of ~70 individuals
552 per species from Southern France natural stands. The SNPs were chosen so that the *illuSNP*
553 set excluded SNPs showing compressed clusters (*i.e.* potential paralogs) and those showing a
554 high number of inconsistencies with control genotypes, as recommended by the authors.
555 Comparing between datasets, for SNPs exclusive to one species in the *sanSNP* set, more than
556 68% either show the same pattern in the *illuSNP* set, or one where the alternative allele was at
557 a frequency below 5% in the other species. Less than 8% of those SNPs are common in both
558 species in the *illuSNP* set. Similarly, for singletons in the *sanSNP* set, more than two-third of
559 the corresponding SNPs in the *illuSNP* set showed very low to low frequency (<10%), while
560 only 11% in *Q. petraea* and 9% in *Q. robur* showed a *maf* above 0.25. This further confirms

561 the reality of singletons in our *DiP*, and also that some may represent more frequent
562 polymorphisms in larger samples of local populations. The correlations among *maf* in both
563 datasets were high and significant (0.66 and 0.68 respectively for *Q. petraea* and *Q. robur*,
564 both $Pr < 0.0001$).

565 Finally, various methodological steps and obtained results tend to demonstrate that we
566 avoided a bias towards low-diversity genic regions: (i) an initial verification that very low
567 BlastX *E*-values ($< 10^{-80}$) did not target more conserved regions, (ii) a primer design
568 optimizing the amplification of polymorphic fragments, both (i) and (ii) using potential
569 variants in ESTs data assembled across both species (Fig. S2-B steps 1 and 3; Appendix S1,
570 Supporting information), (iii) a high nucleotide diversity across genes and ~50% of shared
571 variants (Table 3 and Fig. 4), (iv) a very low proportion of fragments with no detected
572 variants, and a substantial part (~30%) of variant positions due to Indels and SSRs (Table 2),
573 (v) additional results showing that, across ~100 kb of more than 150 independent fragments
574 amplifying in one species only and thus with possible more divergent primer pairs, the
575 number of detected heterozygotes was twice smaller compared to fragments amplifying in
576 both species (more details in Appendix S1, Supporting information).

577 These results altogether suggest a small risk of SNP ascertainment bias if these new resources
578 were to be used in populations both within and/or outside the geographic distribution
579 surveyed, in contrast to panels with much less individuals than here (see respectively
580 Lepoittevin *et al.* 2015 for a discussion on the consequences of such bias in *Quercus* species,
581 and Temunovic *et al.* 2020 cited above).

582 Overall, we obtained sequence data for 0.072% (~530 kb) of the haploid genome of *Q. robur*
583 (size of ~740 Mb in Kremer *et al.* 2007). We also targeted ~3% of the 25808 gene models
584 described in the oak genome sequencing project (www.oakgenome.fr), and around 1% of the
585 gene space in length. Interestingly, both randomly chosen genic regions and those covering
586 different functional categories have been mapped across all linkage groups (columns F and X
587 in Table S1, Supporting information). Due to the absence of observed background LD, their
588 diversity patterns can be considered independent. The genes studied represent a large number
589 of categories, as illustrated by very similar distributions for level 2 GO terms to those
590 obtained with the larger *ocv4* assembly (Lesur *et al.* 2015, comparing their Figure 2 to Fig.
591 S4-A to S4-C, Supporting information).

592 *Diversity magnitude and heterogeneity highlight species integrity and introgression patterns*

593 Using a detailed polymorphism typology, we characterized for the first time in two oak
594 species a high proportion of variant positions (30%) that included 1 bp to medium-sized
595 indels and sequence repeats, compared to the more common and commonly reported SNP loci
596 (Table 2). The proportions of indels observed (11.5% of all polymorphisms) is in the range of
597 results available in model tree species (e.g. 13.8% across the genome in *Prunus avium*,
598 Shirasawa *et al.* 2017; 19% in *Prunus persica*, Cao *et al.* 2014; a lower estimate of 1.4% in
599 *Populus trichocarpa*, Evans *et al.* 2014). Although less abundant than SNPs, they represent an
600 important component of nucleotide variation, often having high functional impacts when
601 located within coding sequences, and they have been proposed as an easy source of markers
602 for natural populations studies (Väli *et al.* 2008). Larger-sized indels are also likely to be
603 relatively frequent in intergenic regions of the *Quercus* genome and have been linked to
604 transposable elements (TE, see the BAC clones overlapping regions analyses in Plomion *et al.*
605 2016). Similarly, large indels and copy number variation linked to TE activity were identified
606 as an important component of variation among hybridizing *Populus* species (Pinosio *et al.*
607 2017). Here when considering variant positions involved in complex polymorphisms, we
608 observed one variant position per 48 bp on average within species (resp. one per 30 bp in
609 both), compared to the one SNP per 68 bp statistic (resp. one SNP per 42 bp across both
610 species). Also, some of the SNPs observed were located within complex polymorphic regions
611 that would have been classically filtered out, and nucleotide diversity (π) estimates were
612 higher by 12% when including all polymorphisms (from 0.0038 to 0.0044 if averaging across
613 both species and all genes, Table 3). These nucleotide diversity estimates are provided for the
614 first time in *Q. petraea* and *Q. robur* across a large number genic regions (> 850), compared
615 to previous candidate genes studies across much smaller numbers (< 10) of gene fragments
616 (Kremer *et al.* 2012 in *Q. petraea*; e.g. Homolka *et al.* 2013).

617 Based on these data, there is an interest in attempting to estimate SNP numbers across the full
618 genome of the studied species for range-wide samples, as it may impact filtering strategies in
619 pipelines for future NGS haplotype-based data production, or decisions to develop or not SNP
620 arrays in these species. In order to do that, a few realistic assumptions can be made from both
621 the exhaustive description of variants provided, and the mean proportions of SNP numbers in
622 new individuals that we computed for increasing across sample sizes. First ~10% additional
623 rare SNPs per sample could be observed for a *DiP* twice as large as ours (based on Fig. S7-A
624 data, Supporting information). Thus given the representativity of our data compared to the

625 *ocv4* unigene (Lesur *et al.* 2015), we would expect around 1.36 million SNPs on average
626 within species by applying our statistics to the full genic partition of *Q. robur* or *Q. petraea*
627 (~80 Mb, www.oakgenome.fr, Plomion *et al.* 2018). Another reasonable assumption is that
628 shared and exclusive polymorphisms proportions across genic regions would be around 30%
629 and 70% respectively, for these closely related oak species (based on both our *DiP* and
630 Lepoittevin *et al.* 2015 results), which translates into the presence of ~2.32 million SNPs for
631 the genic partition in a sample including both *Q. petraea* and *Q. robur* (resp. ~4.22 if
632 including also *Q. pubescens* and *Q. pyrenaica*). Finally, if we apply to the *Quercus* genome a
633 range of ratios for SNPs counts in intergenic over genic regions estimated from several tree
634 species natural population samples (2.03 in *Populus trichocarpa*, Zhou and Holliday 2012;
635 2.25 in the “3P” *Q. robur* reference genotype, Plomion *et al.* 2016; 2.57 in *Prunus persica*
636 wild accessions, Cao *et al.* 2014), we obtain an estimate of between 34 to 42 million SNPs
637 within species across a large spatial range (resp. 41 to 51 million SNPs in both *Q. petraea* and
638 *robur* species, and 75 to 94 million SNPs considering the 4 species previously cited). All these
639 figures could be at least 30% higher if one considers all possible variants involved in indels,
640 SSRs and complex polymorphisms, as shown in our results. Although of the same order of
641 magnitude, the contrast with the twice smaller number of SNPs identified in Leroy *et al.* 2019
642 (~32 millions) across the same four species with similar sample sizes than ours, could be
643 explained by different factors. First their filtering strategy applied on Pool-seq data in order to
644 minimize errors basically excludes all singletons. However, we have seen that verified
645 singletons which could represent rare or local variants amounted to more than 20% of all
646 polymorphisms (see Results). Indeed, very stringent filters are often applied in practice to
647 limit error rates and avoid false-positives, hence limiting the impact of variable read depth and
648 possible ascertainment bias risks, which altogether significantly decrease the number of
649 informative loci compared to either initial fixed amounts (in genotyping arrays, e.g.
650 Lepoittevin *et al.* 2015) or potential amounts (in reference genomes, e.g. Pina-Martins *et al.*
651 2019 in *Quercus* species; see also Van Dijk *et al.* 2014). Second, no cross-validation step is
652 available in Leroy *et al.* (2019) for data quality, that would have permitted to have a better
653 grasp of possible bias and error rate expected in such a dataset, and its consequences on allele
654 frequency estimates and inference methods (see Hivert *et al.* 2018 and discussion below).
655 Also, we can't exclude that a regional sampling strategy such as the one used in Leroy *et al.*
656 (2019) might miss allelic variants with a higher *maf* in other regions for the two species
657 having the wider geographical range.

658 Our nucleotide diversity estimates are consistent with those obtained from genome-wide data
659 and range-wide panels in angiosperm tree species, available mostly from the model genus
660 *Populus* (e.g. *P. trichocarpa*: 1 SNP per 52 bp and $\pi \sim 0.003$ across genic regions, Zhou and
661 Holliday 2012, Zhou *et al.* 2014, Evans *et al.* 2014, Wang *et al.* 2016; *P. tremula*: $\pi \sim 0.008$, *P.*
662 *tremuloides*: $\pi \sim 0.009$ across genic regions, Wang *et al.* 2016; $\pi \sim 0.0026$ to 0.0045 in a panel
663 including wild *Prunus persica* accessions, Cao *et al.* 2014). These diversity levels are also
664 within the range estimated for the long-term perennial outcrosser category in Chen *et al.*
665 (2017, see Fig. 1-D with a mean value of silent π close to ~ 0.005) and can be considered
666 relatively high in the plant kingdom if excluding annual outcrosser estimates or intermediate
667 otherwise. In oaks as in many other tree species with similar life history traits, these high
668 levels would be consistent with their longevity, large variance in reproductive success and
669 recolonization or introgression histories, which could have maintained deleterious loads of
670 various origins (Zhang *et al.* 2016, Chen *et al.* 2017, Christe *et al.* 2016b).

671 Comparing the nucleotide diversity distributions and examining the range of differentiation
672 across genic regions in our *Dip* reveal several robust patterns that altogether illustrate
673 historical introgression among both *Quercus* species. These two species have long been
674 considered as iconic examples of species exhibiting high levels of gene flow (e.g: Petit *et al.*
675 2003; Arnold 2006), despite more recent evidence of strong reproductive barriers (Abadie *et*
676 *al.* 2012). What has been referred to as “strong species integration” seems nevertheless clearer
677 in our *Dip* for *Q. robur* than for *Q. petraea*, according to genetic clustering inference without
678 any *a priori*. Three individuals (27%) considered as typical morphological *Q. petraea* adults
679 (Kremer *et al.* 2002a) showed significant levels of introgression (Fig. 3). In contrast, only one
680 *Q. robur* based on morphology was introgressed to a level matching the least introgressed *Q.*
681 *petraea* individual. Discussing species delimitation, Guichoux *et al.* (2013) also showed more
682 robustness in assigning morphological *Q. robur* individuals to their genetic cluster,
683 illustrating an asymmetry in their introgression levels. We note that among our *Dip*
684 individuals, Qs28, one parent from two mapping pedigrees (Bodénès *et al.* 2016) is a clear F1
685 hybrid among both species (Fig. 3), making those pedigrees two back-crosses instead of one
686 cross within species and one between species.

687 Moreover, after excluding the four most introgressed individuals, nucleotide diversity in *Q.*
688 *petraea* was significantly higher (by $\sim 5\%$ on average) than in *Q. robur*. This effect is small,
689 detectable only with Wilcoxon paired ranked tests, mostly across the same ~ 200 regions
690 sampled randomly and in the *Abiotic stress* category, despite the very large diversity variance

691 across regions, and robust to excluding the highest diversity values. We also sequentially
692 removed the three individuals with the highest Q -values from the *Q. petraea* cluster (Fig. 3),
693 since they could still harbor residual heterozygosity due to recent back-crossing events and
694 generate the pattern observed. Remarkably, the same significant patterns of higher diversity in
695 *Q. petraea* were observed. Therefore, with 8 to 10 gametes in *Q. petraea* instead of 8 to 24
696 gametes in *Q. robur*, and with twice less natural stands sampled, the nucleotide diversity in *Q.*
697 *petraea* was still slightly and significantly higher than in *Q. robur* ($Pr < 0.011$ and $Pr < 0.026$,
698 using all polymorphisms or SNPs only respectively). Although the magnitudes of range-wide
699 population structure within both species could differentially affect both species global
700 diversity across our *Dip*, published results show that these are very small with similar values
701 ($\sim 1\%$ across SNPs, Guichoux *et al.* 2013).

702 The main hypotheses proposed so far to explain this difference in extent of diversity between
703 species relate to their disparities in life-history strategies for colonizing new stands and
704 associated predictions (Petit *et al.* 2003, Guichoux *et al.* 2013). The colonization dynamics
705 model and patterns observed also assumes very similar effective population sizes in both
706 species, which is a reasonable assumption due to their shared past history and the strong
707 introgression impact at the genomic level. However, given increasing and recent evidence of
708 pervasive effects of different types of selection across genic regions with high-throughput
709 data (e.g. Zhang *et al.* 2016; Christe *et al.* 2016b in *Populus*; Chen *et al.* 2017 for long-term
710 perennials), alternative (and non-exclusive) hypotheses worth considering are ones of a higher
711 genome-wide impact of selective constraints in *Q. robur* (Gillespie 2000; Hahn 2008; Cutter
712 and Payseur 2013; Kern and Hahn 2018; e.g. Grivet *et al.* 2017). Since *Q. robur* is the most
713 pioneering species, it has likely been submitted to very strong environmental pressures at the
714 time of stand establishment. Selection might be efficient, given oak tree reproductive
715 capacities, and affect variation across a large number of genes involved in abiotic and biotic
716 responses. This would be consistent with significantly lower levels of diversity (He) in *Q.*
717 *robur* at SNPs located in genes that were specifically enriched for abiotic stress GO terms
718 (Guichoux *et al.* 2013, see their Table S5). Redoing here the same tests across a larger number
719 of independent SNPs (> 1000), *Q. petraea* systematically showed the same trend of a slightly
720 higher diversity overall, and significantly so only for the *Abiotic stress* category ($Pr < 0.01$)
721 and for a similar outlier SNP category ($F_{ST} > 0.4$, mean $He > 0.15$, $Pr < 0.001$) than in Guichoux
722 *et al.* (2013). In summary, the absence of the same pattern in any other functional categories
723 might suggest that these are too broad in terms of corresponding biological pathways, hence

724 mixing possible selection signals of opposite effects among species, while we still detect an
725 overall effect due to linked selection on a random set of genes, and on genes involved in
726 abiotic stress.

727 Within both species, no differences in nucleotide diversity, and a very small differentiation
728 (below 1.5%) were found on average across genes among the main cpDNA lineages (B *versus*
729 A or C) that indicate past refugial areas and migration routes. These patterns were expected,
730 given oaks' life history traits (e.g. high fecundity and dispersal rates), large population sizes,
731 and plausible recolonization scenarios throughout Europe leading to current adaptive
732 differentiation among populations at both nuclear genes and traits (Kremer *et al.* 2010). Only
733 cpDNA ancient differentiation signals among isolated historical refugia were retained, while
734 other putative adaptive divergence effects due to different environments were erased, as
735 illustrated by an absence of correlations between cpDNA and nuclear or phenotypic traits
736 divergence across populations (Kremer *et al.* 2002b). This is consistent with many events of
737 population admixture during the last ~6000 thousands years after European regions were
738 recolonized, as well as a very low genetic differentiation among distant populations (e.g.
739 Guichoux *et al.* 2013), which contrasts with a much higher differentiation often observed for
740 adaptive traits (e.g. Kremer *et al.* 2014; Sáenz-Romero *et al.* 2017). Interestingly, the very
741 few genes with mean F_{ST} between 0.21 and 0.56 among lineages are not the same in *Q.*
742 *petraea* and *Q. robur* (five and seven genes respectively). Seven of them have GO terms
743 indicating their likely expression in chloroplasts, or their interaction with chloroplastic
744 functions. They are either housekeeping genes for basic cellular functions, or belong to biotic
745 or abiotic stress functions (seven of them), and could be involved in local adaptation between
746 ecologically distant populations, calling for further research in larger samples.

747 More generally, analyses comparing the nucleotide diversity patterns at genes involved in
748 both species relevant biosynthesis pathways for ecological preferences (e.g. Porth *et al.* 2005;
749 Le Provost *et al.* 2012, 2016) are clearly needed in replicated populations, for example to
750 estimate the distribution and direction of selection effects and putative fitness impact across
751 polymorphic sites (Stoletzki and EyreWalker 2011), or to study the interplay between
752 different types of selection and variation in local recombination rates on both diversity and
753 differentiation patterns (Payseur and Rieseberg 2016).

754 A large proportion of shared polymorphic sites (~50% in any species) highlights the close
755 proximity of species at the genomic level, consistently with a low mean differentiation across
756 polymorphic sites (F_{ST} ~0.13, Fig. 4-C), and despite the very large heterogeneity observed

757 across differentiation estimates. This has now been classically interpreted (and modeled) as
758 reflecting a strong variance in migration and introgression rates, in oaks in particular (Leroy *et*
759 *al.* 2017), with islands of differentiation assumed to represent regions resistant to
760 introgression. However, interpretations of such patterns remain controversial and multiple
761 processes might be involved and worth exploring further in oaks, such as the effects of
762 heterogeneous selection (both positive and background) at linked loci (Cruickshank and Hahn
763 2014; Wolf and Ellegren 2017). These effects could be particularly visible in low-
764 recombination regions (Ortiz-Barrientos *et al.* 2016), and would further interact with the
765 mutational and recombination landscapes during the course of speciation (Ortiz-Barrientos
766 and James 2017) and during their complex demographic history.

767 *Applications and usefulness as reference data*

768 During this project, several studies valued part of these resources, hence illustrating their
769 usefulness. For example, good quality homologous sequences were also obtained for ~50 %
770 of the gene fragments in one individual of *Quercus ilex*. This species is relatively distant
771 genetically to both *Q. petraea* and *Q. robur*, belonging to a different section, so these data
772 guided the choice of nuclear genes for better inferring phylogenetic relationships across 108
773 oak species (Hubert *et al.* 2014). Bioinformatics tools and candidate genes annotated during
774 the project were also useful to similar genes and SNP discovery approach in *Quercus* or more
775 distant *Fagaceae* species (Rellstab *et al.* 2016, Lalagüe *et al.* 2014 in *Fagus sylvatica*, El
776 Mujtar *et al.* 2014 in *Nothofagus* species). Given the low ascertainment bias and good
777 conversion rate expected within the range surveyed, those genomic resources would be
778 directly applicable to landscape genomics studies at various spatial scales (reviewed in Fetter
779 *et al.* 2017) in both *Quercus* species. Indeed, easy filtering on provided SNP statistics in the
780 catalog would allow distinguishing among different classes of SNPs (e.g. exclusive to each
781 species, common and shared by both, linked to particular GO functional categories),
782 delimiting and tracing species in parentage analyses and conservation studies (e.g. Guichoux
783 *et al.* 2013; Blanc-Jolivet *et al.* 2015), or improving estimates of lifetime reproductive success
784 and aiming to understand how demographic history and ecological drivers of selection affect
785 spatial patterns of diversity or isolating barriers (Andrew *et al.* 2013; e.g. Geraldès *et al.*
786 2014). This type of spatial studies are surprisingly rare in these oak species, they usually
787 include a small number of SSR markers, and all suggest complexity in geographical patterns
788 of genetic variation and importance of the ecological context (e.g. Neophytou *et al.* 2010;
789 Lagache *et al.* 2014; Klein *et al.* 2017, Beatty *et al.* 2016 for local or regional studies; Muir

790 and Schlötterer 2005; Gerber *et al.* 2014, Porth *et al.* 2016 for range-wide studies). Their
791 power and scope would likely be greatly improved by using medium-scale genotyping dataset
792 including a few thousands SNPs such as those described in our study.

793 The robust patterns described above of differentiation heterogeneity and consistent
794 differences in diversity magnitude among species call for more studies at both spatial and
795 genomic scales for unraveling these species evolutionary history, in particular regarding the
796 timing, tempo, dynamics and genetic basis of divergence and introgression. Practically, in
797 order to address those questions in oaks, genomic data on larger samples of individuals could
798 be obtained from either genome complexity reduction methods such as RAD-seq and similar
799 approaches (e.g. Andrews *et al.* 2016) or previously developed SNP arrays (e.g. Silva-Junior
800 *et al.* 2015). We do not recommend the development of a very large SNP array in oaks since it
801 is likely to be very costly for the actual return, especially given the very large and range-wide
802 panel that would be needed to significantly limit ascertainment bias (see Lepoittevin *et al.*
803 2015). The very low overall levels of LD observed here indicate also potentially high
804 recombination rates, and thus that a very high SNP density would be required for targeting
805 functional variants, which would not be compatible with technical constraints for controlling
806 for genotyping error rates (previously shown to be high in SNP array). Indeed, these rates
807 would probably be stronger for high diversity, complex, duplicate or multiple copy genic
808 regions (as those observed in this study in Tables S1 and S4, Supporting information, and
809 shown recently to have an evolutionary impact on the *Q. robur* genome structure, Plomion *et al.*
810 2018), preventing these regions to be included in SNP arrays. The very short LD blocks
811 observed in this study might also limit the utility of RADseq data alone to uncover many loci
812 potentially under selection in genome scans for local adaptation studies (Lowry *et al.* 2016;
813 McKinney *et al.* 2017). In contrast, targeted sequence capture (TSC) strategies for
814 resequencing (Jones and Good 2016), and the more recent advances in RADseq approaches
815 that deal with previous limitations (Arnold *et al.* 2013; Henning *et al.* 2014; and see Rochette
816 *et al.* 2019), although still uncommon in forest tree species evolutionary studies, might be
817 more useful and efficient since they can be oriented towards recovering long genomic
818 fragments. They would thus allow more powerful site frequency spectrum and haplotype-
819 based inferences to be pursued, therefore avoiding most of the SNP arrays technical issues
820 (e.g. Zhou *et al.* 2014; Wang *et al.* 2016), especially given the large variance in nucleotide
821 diversity and low overall differentiation characterized here. TSC approaches will surely be
822 encouraged and tailored to specific evolutionary research questions in oaks in the next decade,

823 given the new *Q. robur* genome sequence availability (Plomion *et al.* 2018; Lesur *et al.* 2018
824 for the first TSC in oaks). However, the bioinformatics pipelines needed for validating
825 haplotype-based or quality data for population genetics inferences also need constant
826 reassessment according to research questions and chosen technology.

827 We thus propose, in addition to direct applications to landscape genetics (detailed above) and
828 transferability to other *Quercus* species (for example using primer information in Table S1,
829 Supporting information, and see Chen *et al.* 2016), that the high-quality data characterized in
830 this study serve as a reference for such validation purposes. They could not only help for
831 adjusting parameters in pipelines for data outputs, but also allow estimating genotyping error
832 rates for SNP and more complex classes of variants, either by comparing general patterns (e.g
833 *maf* distribution from Tables S3, S4 Supporting information) or using the same control
834 individuals maintained in common garden that could be included in larger-scale studies. Such
835 a reference catalog of SNPs and other types of polymorphisms within gene fragments could
836 also be very useful for solid cross-validation of variants identification, allele frequency and
837 other derived summary statistics in alternative strategies such as *Pool-Seq*, which allow
838 increasing genomic coverage while sampling cost-effectively by pooling individuals
839 (Schlötterer *et al.* 2014). Indeed, the drawback of *Pool-Seq* approaches, despite dedicated
840 software (PoPoolation2, Kofler *et al.* 2011) is that they can give strongly biased estimates, or
841 ones that do not consider evolutionary sampling (Hivert *et al.* 2018). Therefore, they require
842 further validation methods which usually value previously developed high-quality and lower-
843 scale data (e.g. *Pool-Seq versus Sanger* and *Rad-Seq* in Christe *et al.* 2016b; *Illumina GA2*
844 *versus Sanger* in Cao *et al.* 2014; *EUChip60K versus deep-whole genome resequencing* in
845 Silva-Junior *et al.* 2015). Finally such a reference dataset would help optimizing the amount
846 of data recovery from either TSC or whole-genome resequencing experiments in future
847 research challenges by fine-tuning dedicated data processing bioinformatics pipelines.

848 **Data Accessibility**

849 The original assembly used for selecting contigs is in Appendix S2 (Supporting information).
850 For Sanger trace files (with data on at least 2 individuals), see the Dryad repository (at the
851 <https://doi.org/10.5061/dryad.4mw6m906j> link). Consensus sequences are respectively in
852 appendices S3 (used to design primers and for functional annotation, see also Table S2), S4
853 (genomic sequences obtained), and S5 (genomic sequences obtained for *Q. ilex*). Tables S1
854 and S2 correct and extend the oak Candidate Genes Database of the Quercus Portal
855 (www.evoltree.eu/index.php/e-recources/databases/candidate-genes). SNP, indel and SSR

856 catalogs and positions within genomic consensus sequences, and ready-to-use format for
857 genotyping essays are provided in Tables S3 to S5 (Supporting information), and at
858 <https://github.com/garniergere/Reference.Db.SNPs.Quercus> with additional information.

859 Bioperl scripts from the SeqQual pipeline are given at
860 <https://github.com/garniergere/SeqQual>, example of parameter files and scripts for
861 STRUCTURE analyses and parsing MREPS software are given at
862 <https://github.com/garniergere/Reference.Db.SNPs.Quercus>

863

864 **Acknowledgments**

865 The authors thank Alexis Ducouso, Jean-Marc Louvet, Guy Roussel, Pablo Goicoechea,
866 Hervé le Bouler, Félix Gugerli, Csaba Matyas, Sandor Bordacs, Hans P. Koelewijn, Joukje
867 Buiteveld, Stephen Cavers, Bernd Degen and Jutta Buschbom for choosing trees and
868 providing dried leaves of individuals from various Intensive Study Populations of previous
869 European projects populations. We are grateful to H. Lalagüe, G. Vendramin, I. Scotti, and L.
870 Brousseau for testing earlier scripts of SeqQual and to I. Lesur for help in using the *ocv4* oak
871 resources. The sequencing work was funded by the EVOLTREE network of Excellence (EU
872 contract n°016322). TL post-doc fellowship was funded by the ANR TRANSBIODIV (06-
873 BDIV-003-04) and LINKTREE (contract n°2008-966). TD salary was funded by the ANR
874 REALTIME (N°59000256). Computing facilities of the Mésocentre de calcul Intensif
875 Aquitain des Universités de Bordeaux, de Pau et des Pays de l'Adour are thanked for
876 providing computer time for this study. We also thank Rémy Petit for funding part of TL
877 fellowship and support in developing SeqQual tools. PA received a Ph.D. grant (2009-2011)
878 from the « Ministère de l'Education Nationale, de l'Enseignement Supérieur et de la
879 Recherche » of France, and additional funding from EVOLTREE. We thank Oliver Brendel,
880 Ricardo Alia, Komlan Avia and Hilke Schröder for reviewing the manuscript and for their
881 constructive comments.

882 **Conflict of interest disclosure**

883 The authors of this article declare that they have no financial conflict of interest with the
884 content of this article.

885 **References**

886 Abadie P, Roussel G, Dencausse B, *et al.* (2012) Strength, diversity and plasticity of
887 postmating reproductive barriers between two hybridizing oak species (*Quercus robur* L. and
888 *Quercus petraea* (Matt) Liebl.). *Journal of Evolutionary Biology*, **25**, 157-173.

- 889 Abbott RJ, James JK, Milne RI, Gillies ACM (2003) Plant introductions, hybridization and
890 gene flow. *Philosophical Transactions of the Royal Society of London B: Biological Sciences*,
891 **358**, 1123–1132.
- 892 Andrew RL, Bernatchez L, Bonin A *et al.* (2013) A road map for molecular ecology.
893 *Molecular Ecology*, **22**, 2605–2626.
- 894 Andrews KR, Good JM, Miller MR, Luikart G, Hohenlohe PA (2016) Harnessing the power
895 of RADseq for ecological and evolutionary genomics. *Nature Reviews Genetics*, 17(2): 81-
896 92.
- 897 Arnold ML (2006) Evolution through genetic exchange. Oxford University Press, Oxford.
- 898 Arnold B, Corbett-Detig RB, Hartl D, Bomblies K (2013) RADseq underestimates diversity
899 and introduces genealogical biases due to nonrandom haplotype sampling. *Molecular biology*
900 **22**: 3179-3190.
- 901 Beatty GE, Montgomery WI, Spaans F, Tosh DG, Provan J (2016) Pure species in a
902 continuum of genetic and morphological variation: sympatric oaks at the edge of their range.
903 *Annals of Botany*, **117**, 541-549.
- 904 Blanc-Jolivet C, Liesebach M (2015) Tracing the origin and species identity of *Quercus robur*
905 and *Quercus petraea* in Europe: a review. *Silvae Genetica* **64(4)**, 182–193.
- 906 Bodénès C, Labbe T, Pradère S, Kremer A (1997) General vs. local differentiation between two closely related
907 white oak species. *Molecular Ecology*, **6**: 713-724.
- 908 Bodénès C, Chancerel E, Ehrenmann F, Kremer A, Plomion C (2016) High-density linkage
909 mapping and distribution of segregation distortion regions in the oak genome. *DNA Research*,
910 **23**, 115-124.
- 911 Bodénès C, Chancerel E, Gailing O, *et al.* (2012) Comparative mapping in the Fagaceae and
912 beyond with EST-SSRs. *BMC Plant Biology*, **12**, 153.
- 913 Bodénès C, Chancerel E, Murat F, *et al.* (2012) Comparative mapping in the Fagaceae and
914 beyond using EST-SSRs. *BMC Plant Biology*, **12**, 153.
- 915 Brewer S, Cheddadi R, De Beaulieu JL, Reille M, Data contributors (2002) The spread of
916 deciduous *Quercus* throughout Europe since the last glacial period. *Forest Ecology and*
917 *Management*, **156**, 27–48.
- 918 Brousseau L, Tinaut A, Duret C, *et al.* (2014) High-throughput transcriptome sequencing and
919 preliminary functional analysis in four neotropical tree species. *BMC Genomics*, **15**, 238.
- 920 Buerkle CA, Gompert Z, Parchman TL (2011) The n=1 constraint in population genomics.
921 *Molecular Ecology*, **20**, 1575–1581.

- 922 Cannon CH, Brendel O, Deng M *et al.* (2018) Gaining a global perspective on *Fagaceae*
923 genomic diversification and adaptation. *New Phytologist*, **218**, 894-897.
- 924 Cao K, Zheng Z, Wang L *et al.* (2014) Comparative population genomics reveals the
925 domestication history of the peach, *Prunus persica*, and human influences on perennial fruit
926 crops. *Genome Biology*, **15**, 415.
- 927 Casasoli M, Derory J, Morera-Dutrey C, *et al.* (2006) Comparison of QTLs for adaptive traits
928 between oak and chestnut based on an EST consensus map. *Genetics*, **172**, 533–546.
- 929 Cavender-Bares J (2016) Diversity, distributions, and ecosystem services of the North-
930 American oaks. *International oaks*, **27**, 37-48.
- 931 Chen J, Glémin S, Lascoux M (2017) Genetic diversity and the efficacy of purifying selection
932 across plant and animal species. *Molecular Biology and Evolution*, **34**, 1417–1428.
- 933 Chen J, Zeng Y-F, Liao W-J *et al.* (2016) A novel set of single-copy nuclear gene markers in
934 white oak and implications for species delimitation. *Tree Genetics & Genomes*, **13**, 50.
- 935 Christe C, Stölting KN, Bresadola L, *et al.* (2016a) Selection against recombinant hybrids
936 maintains reproductive isolation in hybridizing *Populus* species despite F1 fertility and
937 recurrent gene flow. *Molecular Ecology*, **25**, 2482–2498.
- 938 Christe C, Stölting KN, Paris M, *et al.* (2016b) Adaptive evolution and segregating load
939 contribute to the genomic landscape of divergence in two tree species connected by episodic
940 gene flow. *Molecular Ecology*, **26**, 59-76.
- 941 Conesa A, Götz S, Garcia-Gomez JM, *et al.* (2005) Blast2GO: a universal tool for annotation,
942 visualization and analysis in functional genomics research. *Bioinformatics*, **21(18)**, 3674–
943 3676.
- 944 Cruickshank TE, Hahn MW (2014) Reanalysis suggests that genomic islands of speciation are
945 due to reduced diversity, not reduced gene flow. *Molecular Ecology*, **23**, 3133–3157.
- 946 Currat M, Ruedi M, Petit RJ, Excoffier L (2008) The hidden side of invasions: massive
947 introgression by local genes. *Evolution*, **62**, 1908–1920.
- 948 Curtu AL, Gailing O, Finkeldey R (2007) Evidence for hybridization and introgression within
949 a species-rich oak (*Quercus* spp.) community. *BMC Evolutionary Biology*, **7**, 218.
- 950 Cutter AD, Payseur BA (2013) Genomic signatures of selection at linked sites: unifying the
951 disparity among species. *Nature Reviews Genetics*, **14**, 262–274.
- 952 Derory J, Scotti-Saintagne C, Bertocchi E, *et al.* (2010) Contrasting relationships between the
953 diversity of candidate genes and variation of bud burst in natural and segregating populations
954 of European oaks. *Heredity*, **104**, 438-448.

- 955 Durand J, Bodénès C, Chancerel E, *et al.* (2010) A fast and cost-effective approach to develop
956 and map EST-SSR markers: oak as a case study. *BMC Genomics*, **11**, 570.
- 957 El Mujtar VA, Gallo LA, Lang T, Garnier-Gere P (2014) Development of genomic resources
958 for *Nothofagus* species using next-generation sequencing data. *Molecular Ecology Resources*,
959 **14**, 1281–1295.
- 960 Evans LM, Slavov GT, Rodgers-Melnick E, *et al.* (2014) Population genomics of *Populus*
961 *trichocarpa* identifies signatures of selection and adaptive trait associations. *Nature Genetics*,
962 **46**, 1089–1096
- 963 Evans LM, Slavov GT, Rodgers-Melnick E, *et al.* (2014) Population genomics of *Populus*
964 *Trichocarpa* identifies signatures of selection and adaptive trait associations. *Nature Genetics*,
965 **46**, 1089-1096.
- 966 Eveno E, Collada C, Guevara MA, *et al.* (2008) Contrasting patterns of selection at *Pinus*
967 *pinaster* Ait. drought stress candidate genes as revealed by genetic differentiation analyses.
968 *Molecular Biology and Evolution* **25**: 417-437.
- 969 Ewing B, Hillier L, Wendl MC, Green P (1998) Base-calling of automated sequencer traces
970 using phred. I. Accuracy assessment. *Genome research*, **8**, 175–185.
- 971 Excoffier L, Lischer HEL (2010) Arlequin suite ver 3.5: a new series of programs to perform
972 population genetics analyses under Linux and Windows. *Molecular Ecology Resources*, **10**,
973 564–567.
- 974 Excoffier L (2007) Analysis of population subdivision. Pages 980-1020 in Handbook of
975 Statistical Genetics. 3rd ed. DJ Balding, M. Bishop, and C. Cannings, ed. Wiley, Chichester,
976 West Sussex, UK.
- 977 Faivre-Rampant P, Lesur I, Boussardon C *et al.* (2011) Analysis of BAC end sequences in
978 oak, a keystone forest tree species, providing insight into the composition of its genome. *BMC*
979 *Genomics*, **12**, 292.
- 980 Falush D, Stephens M, Pritchard JK (2003) Inference of population structure using multilocus
981 genotype data: linked loci and correlated allele frequencies. *Genetics*, **164**, 1567–1587.
- 982 Fetter KC, Gugger PF, Keller SR (2017) Landscape Genomics of Angiosperm Trees: From
983 historic Roots to Discovering New Branches of Adaptive Evolution. In Groover A. and Cronk
984 Q. (eds) *Comparative and evolutionary genomics of angiosperm trees*, *Plant Genetics and*
985 *Genomics: Crops and Models*. New York, Springer.
- 986 Geraldés A, Farzaneh N, Grassa CJ, *et al.* (2014) Landscape genomics of *Populus*
987 *trichocarpa*: the role of hybridization, limited gene flow, and natural selection in shaping
988 patterns of population structure. *Evolution*, **68**, 3260–80.

- 989 Geraldès A, Pang J, Thiessen N, *et al.* (2011) SNP discovery in black cottonwood (*Populus*
990 *trichocarpa*) by population transcriptome resequencing. *Molecular Ecology Resources*, **11**
991 (Suppl. 1), 81–92.
- 992 Gerber S, Chadœuf J, Gugerli F *et al.* (2014) High rates of gene flow by pollen and seed in
993 oak populations across Europe. *PLoS ONE*, **9**, e85130.
- 994 Gilbert KJ, Andrew RL, Bock DG *et al.* (2012) Recommendations for utilizing and reporting
995 population genetic analyses: the reproducibility of genetic clustering using the program
996 STRUCTURE. *Molecular ecology*, **21**, 4925–4930.
- 997 Gillespie JH (2000) Genetic drift in an infinite population: the pseudohitchhiking model.
998 *Genetics*, **155**, 909–919.
- 999 Grivet D, Avia K, Vaattovaara A, Eckert AJ, Neale DB, Savolainen O, Gonzalez-Martinez
1000 SC. 2017. High rate of adaptive evolution in two widespread European pines. *Molecular*
1001 *Ecology*, **26**, 6857–6870.
- 1002 Grivet D, Deguilloux M-F, Petit RJ, Sork VL (2006) Contrasting patterns of historical
1003 colonization in white oaks (*Quercus* spp.) in California and Europe. *Molecular Ecology* **15**,
1004 4085–93.
- 1005 Gugger PF, Cokus SJ, Sork VL (2016) Association of transcriptome-wide sequence variation
1006 with climate gradients in valley oak (*Quercus lobata*). *Tree Genetics and Genomes*, **12**, 15.
- 1007 Guichoux E, Garnier-Géré P, Lagache L *et al.* (2013) Outlier loci highlight the direction of
1008 introgression in oaks. *Molecular Ecology*, **22**, 450–462.
- 1009 Guichoux E, Lagache L, Wagner S, Léger P, Petit RJ (2011) Two highly-validated multiplex
1010 (12-plex and 8-plex) for species delimitation and parentage analysis in oaks (*Quercus* spp.)
1011 *Molecular Ecology Resources*, **11**, 578–585.
- 1012 Hahn MW. 2008. Toward a selection theory of molecular evolution. *Evolution* 62(2):255–
1013 265.
- 1014 Hedrick PW (2005) A standardized genetic differentiation measure. *Evolution*, **59**, 1633–
1015 1638.
- 1016 Henning F, Lee HJ, Franchini P, Meyer A (2014) Genetic mapping of horizontal stripes in
1017 Lake Victoria cichlid fishes: benefits and pitfalls of using RAD marker for dense linkage
1018 mapping. *Molecular Ecology*, **23**, 5224–5240.
- 1019 Hewitt GM (2000) The genetic legacy of the Quaternary ice ages. *Nature* **405**, 907–913.
- 1020 Hivert V, Leblois R, Petit EJ, Gautier M, Vitalis R (2018) Measuring genetic differentiation
1021 from Pool-seq data. *Genetics*, **210**, 315–330.

- 1022 Holliday JA, Aitken SN, Cooke JEK, et al. (2017) *Advances in ecological genomics in forest*
1023 *trees and applications to genetic resources conservation and breeding. Molecular Ecology*, **26**,
1024 706–717.
- 1025 Homolka A, Schueler S, Burg K, Fluch S, Kremer A (2013) Insights into drought adaptation
1026 of two European oak species revealed by nucleotide diversity of candidate genes. *Tree*
1027 *Genetics & Genomes*, **9**, 1179–1192.
- 1028 Hubert F, Grimm GW, Jousselein E, et al. (2014) Multiple nuclear genes stabilize the
1029 phylogenetic backbone of the genus *Quercus*. *Systematics and Biodiversity*, **12**, 405–423.
- 1030 Jeffries DL, Copp GH, Lawson Handley L, et al. (2016) Comparing RADseq and
1031 microsatellites to infer complex phylogeographic patterns, an empirical perspective in the
1032 Crucian carp, *Carassius carassius*, L. *Molecular Ecology*, **25**, 2997–3018.
- 1033 Jensen J, Larsen A, Nielsen LR, Cottrell J (2009) Hybridization between *Quercus robur* and
1034 *Q. Petraea* in a mixed oak stand in Denmark. *Annals of Forest Science*, **66**, 706.
- 1035 Jones MR, Good JM (2016) Targeted capture in evolutionary and ecological genomics.
1036 *Molecular Ecology*, **25**, 185–202.
- 1037 Jones MR, Good JM (2016) Targeted capture in evolutionary and ecological genomics.
1038 *Molecular Ecology*, **25**, 185–202.
- 1039 Keller SR, Olson MS, Silim S et al. (2010) Genomic diversity, population structure, and
1040 migration following rapid range expansion in the Balsam Poplar, *Populus balsamifera*.
1041 *Molecular Ecology*, **19**, 1212–1226.
- 1042 Kern AD, Hahn MW (2018) The neutral theory in light of natural selection. *Molecular*
1043 *Biology and Evolution*, **35**, 1366–1371.
- 1044 Klein EK, Lagache-Navarro L, Petit RJ (2017) Demographic and spatial determinants of
1045 hybridization rate. *Journal of Ecology*, **105**, 29–38.
- 1046 Kofler R, Pandey RV, Schlotterer C (2011) PoPoolation2: identifying differentiation between
1047 populations using sequencing of pooled DNA samples (Pool-Seq). *Bioinformatics*, **27**, 3435–
1048 3436.
- 1049 Kolpakov R, Bana G, Kucherov G (2003) mreps: efficient and flexible detection of tandem
1050 repeats in DNA. *Nucleic Acid Research*, **31**, 3672–3678.
- 1051 Kremer A, Abbott A, Carlson J, et al. (2012) Genomics of Fagaceae. *Tree Genetics &*
1052 *Genomes*, **8**, 583–610.
- 1053 Kremer A, Casasoli M, Barreneche T, et al. (2007) Fagaceae. In: Genome Mapping and
1054 Molecular Breeding in Plants (ed. Kole CR), Vol 7 *Forest Trees*, pp. 165–187. Springer,
1055 Heidelberg, Berlin, New York, Tokyo.

- 1056 Kremer A, Dupouey JL, Deans JD, *et al.* (2002a) Leaf morphological differentiation between
1057 *Quercus robur* and *Quercus petraea* is stable across western European mixed oak stands.
1058 *Annals of Forest Science*, **59**, 777–787.
- 1059 Kremer A, Kleinshmit J, Cotrell J, Cundall EP, Deans JD, *et al.* (2002b) Is there a correlation
1060 between chloroplastic and nuclear divergence, or what are the roles of history and selection on
1061 genetic diversity in European oaks? *Forest Ecology and Management* 156:75-
- 1062 Kremer A, Le Corre V, Petit RJ, Ducouso A (2010) Historical and contemporary dynamics
1063 of adaptive differentiation in European oaks. In *Molecular Approaches in Natural Resource*
1064 *Conservation*. Eds. DeWoody, A., Bickham, J., Michler, C., Nichols, K., Rhodes, G. and
1065 Woeste, K., Cambridge University Press, pp. 101-122.
- 1066 Kremer A, Potts BM, Delzon S (2014) Genetic divergence in forest trees: understanding the
1067 consequences of climate change. *Functional Ecology* **28**, 22–36.
- 1068 Lagache L, Klein EK, Ducouso A, Petit RJ (2014) Distinct male reproductive strategies in
1069 two closely related oak species. *Molecular Ecology*, **23**, 4331–4343.
- 1070 Lalagüe H, Csilléry K, Oddou-Muratorio S, Safrana J, de Quattro C, Fady B, Gonzalez-
1071 Martinez SC, Vendramin GG (2014) Nucleotide diversity and linkage disequilibrium at 58
1072 stress response and phenology candidate genes in a European beech (*Fagus sylvatica*)
1073 population from southeastern France. *Molecular Ecology* **23**, 4696-4708.
- 1074 Lascoux M, Petit RJ (2010) The ‘New Wave’ in plant demographic inference: more loci and
1075 more individuals. *Molecular Ecology*, **19**, 1075–1078.
- 1076 Le Provost G, Lesur I, Lalanne C *et al.* (2016) Implication of the suberin pathway in
1077 adaptation to waterlogging and hypertrophied lenticels formation in pedunculate oak
1078 (*Quercus robur* L.). *Tree Physiology*, **36**, 1330–1342.
- 1079 Le Provost G, Sulmon C, Frigerio JM, *et al.* (2012) Role of waterlogging-responsive genes in
1080 shaping interspecific differentiation between two sympatric oak species. *Tree Physiology*, **32**,
1081 119–134.
- 1082 Lepais O, Gerber S (2011) Reproductive patterns shape introgression dynamics and species
1083 succession within the European white oak species complex. *Evolution*, **65**(1), 156–170.
- 1084 Lepais O, Petit RJ, Guichoux E, *et al.* (2009) Species relative abundance and direction of
1085 introgression in oaks. *Molecular Ecology*, **18**, 2228–2242.
- 1086 Lepais O, Roussel G, Hubert F, Kremer A, Gerber S (2013) Strength and variability of
1087 postmating reproductive isolating barriers between four European white oak species. *Tree*
1088 *Genetics Genomes*, **9**(3), 841–853.

- 1089 Lepoittevin C, Bodénès C, Chancerel E, *et al.* (2015) Single-nucleotide polymorphism
1090 discovery and validation in high density SNP array for genetic analysis in European white
1091 oaks. *Molecular Ecology Resources*, **15**, 1446–1459.
- 1092 Leroy T, Roux C, Villate L, Bodénès C, Romiguier J, Paiva JAP, *et al.* (2017). Extensive
1093 recent secondary contacts between four European white oak species. *New Phytologist*, **214**,
1094 865–878.
- 1095 Leroy T, Rougemont Q, Dupouey J-L, Bodénès C, Lalanne C, Belser C, Labadie K, Le
1096 Provost G, Aury J-M, Kremer A, Plomion C (2019) Massive postglacial gene flow between
1097 European white oaks uncovered genes underlying species barriers. *New Phytologist* early
1098 view, <https://doi.org/10.1111/nph.16039>.
- 1099 Lesur I, Alexandre H, Boury C, *et al.* (2018) Development of target sequence capture and
1100 estimation of genomic relatedness in a mixed oak stand. *Frontiers in Plant Science*
1101 (*METHODS*), doi: 10.3389/fpls.2018.00996.
- 1102 Lesur I, Le Provost G, Bento P, *et al.* (2015) The oak gene expression atlas: insights into
1103 Fagaceae genome evolution and the discovery of genes regulated during bud dormancy
1104 release. *BMC Genomics*, **16**, 112.
- 1105 Lowry DB, Hoban S, Kelley JL *et al.* (2016) Breaking RAD: an evaluation of the utility of
1106 restriction site-associated DNA sequencing for genome scans of adaptation. *Molecular*
1107 *Ecology Resources*, **17**, 142–152.
- 1108 Mariette S, Cottrell J, Csaikl UM, Goikoechea P, Nig A, Lowe AJ, *et al.* (2002) Comparison
1109 of levels of genetic diversity detected with AFLP and microsatellite markers within and
1110 among mixed *Q. petraea* (Matt.) Liebl. and *Q. robur* L. stands. *Silvae Genet.* **51**: 72-79.
- 1111 McKinney GJ, Larson WA, Seeb LW, Seeb JE (2017) RADseq provides unprecedented
1112 insights into molecular ecology and evolutionary genetics: Comment on Breaking RAD by
1113 Lowry et al (2016). *Molecular Ecology Resources* **17**(3), 356–361.
- 1114 Mishra B, Gupta DK, Pfenniger M, *et al.* (2018) A reference genome of the European beech
1115 (*Fagus sylvatica* L.) *GigaScience*, **7**:6. <https://doi.org/10.1093/gigascience/giy063>.
- 1116 Muir G, Fleming CC, Schlötterer C (2000) Species status of hybridizing oaks. *Nature*, **405**,
1117 1016.
- 1118 Muir G, Schlötterer C (2005) Evidence for shared ancestral polymorphism rather than
1119 recurrent gene flow at microsatellite loci differentiating two hybridizing oaks (*Quercus* spp.).
1120 *Molecular Ecology*, **14**, 549–561.

- 1121 Nazareno A. G., Bemmels, J. B., Dick, C. W., & Lohmann, L. G. (2017). Minimum sample
1122 sizes for population genomics: An empirical study from an Amazonian plant species.
1123 *Molecular Ecology Resources*, **17**, 1136–1147.
- 1124 Neale DB, Kremer A (2011) Forest tree genomics: growing resources and applications.
1125 *Nature Reviews Genetics*, **12**, 111–122.
- 1126 Neale DB, Langley CH, Salzberg SL, Wegrzyn JL (2013) Open access to tree genomes: the
1127 path to a better forest. *Genome Biology*, **14**: 120.
- 1128 Neale DB, Martínez-García PJ, La Torre De AR, Montanari S, Wei X-X (2017) Novel in-
1129 sights into tree biology and genome evolution as revealed through genomics. *Annual Reviews*
1130 *of Plant Biology*, **68**, 457–483.
- 1131 Nei M (1987) *Molecular Evolutionary Genetics*. New York, Columbia University Press.
- 1132 Nei M. (1977) F-statistics and analysis of gene diversity in sub-divided populations. *Annals of*
1133 *Human Genetics*, **41**, 225–233.
- 1134 Neophytou C, Gärtner SM, Vargas-Gaete R, Michiels H-G (2015) Genetic variation of
1135 Central European oaks: shaped by evolutionary factors and human intervention? *Tree*
1136 *Genetics & Genomes*, **11**, 79.
- 1137 Nickerson DA, Tobe VO, Taylor SL (1997) PolyPhred: automating the detection and
1138 genotyping of single nucleotide substitutions using fluorescence-based resequencing. *Nucleic*
1139 *Acids Research*, **25**, 2745–2751.
- 1140 Ortiz-Barrientos D, Baack EJ (2014) Species integrity in trees. *Molecular Ecology*, **23**, 4188-
1141 4191.
- 1142 Ortiz-Barrientos D, Engelstädter J, Rieseberg LH (2016) Recombination rate evolution and
1143 the origin of species. *Trends in Ecology and Evolution*, **31**, 226–236.
- 1144 Ortiz-Barrientos D, James ME (2017) Evolution of recombination rates and the genomic
1145 landscape of speciation. *Journal of Evolutionary Biology*, **30**,
- 1146 Parent GJ, Raherison E, Sena J, Mackay JJ (2015) Forest tree genomics: review of progress.
1147 In Plomion C, Adam-Blondonpp A-F (eds) *Land Plants - Trees. Advances in Botanical*
1148 *Research*, **74**, 39–92, London: Academic Press, Elsevier.
- 1149 Payseur BA, Rieseberg LH (2016). A genomic perspective on hybridization and speciation.
1150 *Molecular Ecology*, **25**, 2337– 2360.
- 1151 Petit RJ, Bodénès C, Ducouso A, Roussel G, Kremer A (2003) Hybridization as a
1152 mechanism of invasion in oaks. *New Phytologist*, **161**: 151-164.
- 1153 Petit RJ, Carlson J, Curtu AL, *et al.* (2013) Fagaceae trees as models to integrate ecology,
1154 evolution and genomics. *New Phytologist*, **197**, 369–371.

- 1155 Petit RJ, Brewer S, Bordacs S, *et al.* (2002a) Identification of refugia and post-glacial
1156 colonisation routes of European white oaks based on chloroplast DNA and fossil pollen
1157 evidence. *Forest Ecology and Management* **156**, 49-74.
- 1158 Petit RJ, Csaikl UM, Bordács S, *et al.* (2002b) Chloroplast DNA variation in European white
1159 oaks: phylogeography and patterns of diversity based on data from over 2600 populations.
1160 *Forest Ecology and Management*, **156(1-3)**, 5-26.
- 1161 Pina-Martins JB, Batista J, Pappas G, Paulo OS (2019) New insights into adaptation and
1162 population structure of cork oak using genotyping by sequencing. *Global Change Biology*
1163 **25**:337–350.
- 1164 Pinosio S, Giacomello S, Faivre-Rampant P, *et al.* (2016) Characterization of the poplar pan-
1165 genome by genome-wide identification of structural variation. *Molecular Biology and*
1166 *Evolution*, **33**, 2706–2719.
- 1167 Plomion C, Aury J-M, Amselem J, *et al.* (2016) Decoding the oak genome: public release of
1168 sequence data, assembly, annotation and publication strategies. *Molecular Ecology Resources*,
1169 **16**, 254–265.
- 1170 Plomion C, Aury JM, Amselem J, *et al.* (2018) Oak genome reveals facets of long lifespan.
1171 *Nature Plants*, **4**, 440–452.
- 1172 Porth I, Garnier-Géré P, Klapste J, Scotti-Saintagne, El-Kassaby YA, Burg K, Kremer A
1173 (2016) Species-specific alleles at a β -tubulin gene show significant association with leaf
1174 morphological variation within *Quercus petraea* and *Q. robur* populations. *Tree Genetics &*
1175 *Genomes* **12**: 81.
- 1176 Porth I, Koch M, Berenyi M, *et al.* (2005) Identification of adaptation-specific differences in
1177 mRNA expression of sessile and pedunculate oak based on osmotic-stress induced genes.
1178 *Tree Physiology*, **25**, 1317–1329.
- 1179 Prunier J, Caron S, MacKay J (2017) CNVs into the wild: Screening the genomes of conifer
1180 trees (*Picea* spp.) reveals fewer gene copy number variations in hybrids and links to
1181 adaptation. *BMC Genomics*, **18**, 97.
- 1182 Ramos AM, Usié A, Barbosa P, *et al.* (2018) The draft genome sequence of cork oak.
1183 *Scientific Data*, **5**, 180069, <http://dx.doi.org/10.1038/sdata.2018.69>
- 1184 Rellstab C, Zoller S, Walthert L, *et al.* (2016) Signatures of local adaptation in candidate
1185 genes of oaks (*Quercus* spp.) with respect to present and future climatic conditions. *Molecular*
1186 *Ecology*, **25**, 5907-5924.

- 1187 Rochette NC, Rivera-Colón AG, Catchen JM (2019) Stacks 2: Analytical methods for paired-
1188 end sequencing improve RADseq-based population genomics. *Molecular Ecology*, 28(21),
1189 4737-4754.
- 1190 Rousset F (2008) Genepop'007: a complete re-implementation of the genepop software for
1191 Windows and Linux. *Molecular Ecology Resources*, **8**, 103-106.
- 1192 Roux C, Fraïsse C, Castric V, Vekemans X, Pogson GH, Bierne N (2014) Can we continue to
1193 neglect genomic variation in introgression rates when inferring the history of speciation? A
1194 case study in a *Mytilus* hybrid zone. *Journal of Evolutionary Biology*, 27, 1662-1675.
- 1195 Sáenz-Romero C, Lamy J-B, Ducouso A, Musch B, Ehrenmann F, Delzon S, Cavers S,
1196 Chałupka W, Dağdaş S, Hansen JK *et al.* (2017) Adaptive and plastic responses of *Quercus*
1197 *petraea* populations to climate across Europe. *Global Change Biology* 23: 2831–2847.
- 1198 Sanger F, Nicklen S, Coulson AR (1977) DNA sequencing with chain-terminating inhibitors.
1199 *Proc Natl Acad Sci USA*, 74:5463-5467. Savolainen O, Pyhajarvi T, Knurr T (2007) Gene
1200 flow and local adaptation in trees. *Annual Review of Ecology Evolution and Systematics*, **38**,
1201 595–619.
- 1202 Schlötterer C, Tobler R, Kofler R, Nolte V (2014) Sequencing pools of individuals-mining
1203 genome-wide polymorphism data without big funding. *Nature Reviews Genetics*, 15, 749–
1204 763.
- 1205 Shirasawa K, Isuzugawa K, Ikenaga M, *et al.* (2017) The genome sequence of sweet cherry
1206 (*Prunus avium*) for use in genomics-assisted breeding. *DNA Research*, **24(5)**, 499-508.
- 1207 Silva-Junior OB, Faria DA, Grattapaglia (2015) A flexible multi-species 60K SNP chip
1208 developed from pooled resequencing of 240 *Eucalyptus* tree genomes across 12 species. *New*
1209 *Phytologist* **206**, 1527-1540.
- 1210 Sonah H, Bastien M, Iquira E, *et al.* (2013) An improved genotyping by sequencing (GBS)
1211 approach offering increased versatility and efficiency of SNP discovery and genotyping. *PLoS*
1212 *One* 8(1), e54603.
- 1213 Sork VL, Fitz-Gibbon ST, Puiu D, *et al.* 2016 First draft assembly and annotation of the
1214 genome of a California endemic oak *Quercus lobata* Née (Fagaceae). *G3: Genes Genomes*
1215 *Genetics*, **6 (11)**, 3485–3495.
- 1216 Sousa VC, Peischl S, Excoffier L (2014) Impact of range expansions on current human
1217 genomic diversity. *Current Opinion in Genetics and Development*, **29**, 22-30
- 1218 Stoletzki N, Eyre-Walker A (2011) Estimation of the neutrality index. *Molecular Biology and*
1219 *Evolution*, **28**, 63–70.

- 1220 Tajima F (1983) Evolutionary relationship of DNA sequences in finite populations. *Genetics*,
1221 **105**, 437-460.
- 1222 Tajima F (1993) Measurement of DNA polymorphism. In: *Mechanisms of Molecular*
1223 *Evolution. Introduction to Molecular Paleopopulation Biology*, edited by Takahata, N. and
1224 Clark, A.G., Tokyo, Sunderland, MA: Japan Scientific Societies Press, Sinauer Associates,
1225 Inc., p. 37-59.
- 1226 Tarkka MT, Herrmann S, Wubet T, *et al.* (2013) OakContigDF159.1, a reference library for
1227 studying differential gene expression in *Quercus robur* during controlled biotic interactions:
1228 use for quantitative transcriptomic profiling of oak roots in ectomycorrhizal symbiosis. *New*
1229 *Phytologist*, **199**, 529–540.
- 1230 Tine M, Kuhl H, Gagnaire P-A, *et al.* (2014) European sea bass genome and its variation
1231 provide insights into adaptation to euryhalinity and speciation. *Nature Communications*, **5**,
1232 5770.
- 1233 Tuskan GA, DiFazio S, Jansson S, *et al.* (2006) The genome of black cottonwood, *Populus*
1234 *trichocarpa* (Torr. & Gray). *Science* 2006, **313(5793)**:1596–1604.
- 1235 Ueno S, Le Provost G, Leger V, *et al.* (2010) Bioinformatic analysis of ESTs collected by
1236 Sanger and pyrosequencing methods for a keystone forest tree species: oak. *BMC Genomics*,
1237 **11**, 650.
- 1238 Valbuena-Carabana M, González-Martínez S, Sork V, *et al.* (2005) Gene flow and
1239 hybridisation in a mixed oak forest (*Quercus pyrenaica* Willd. and *Quercus petraea* (Matts.)
1240 Liebl.) in central Spain. *Heredity*, **95**, 457–465.
- 1241 Väli U, Brandström M, Johansson M, Ellegren H (2008) Insertion-deletion polymorphisms
1242 (indels) as genetic markers in natural populations. *BMC Genetics*, **9**, 8.
- 1243 Van Dijk EL, Auger H, Jaszczyszyn Y, Thermes C (2014) Ten years of next-generation
1244 sequencing technology. *Trends in Genetics*, **30**, 418–426.
1245 <https://doi.org/10.1016/j.tig.2014.07.001>
- 1246 Verde I, Abbot GA, Scalabrin S, *et al.* (2013) The high-quality draft genome of peach
1247 (*Prunus persica*) identifies unique patterns of genetic diversity, domestication and genome
1248 evolution. *Nature Genetics*, **45**, 487–494.
- 1249 Verde I, Jenkins J, Dondini L, *et al.* (2017) The Peach v2.0 release: high-resolution linkage
1250 mapping and deep resequencing improve chromosome-scale assembly and contiguity. *BMC*
1251 *Genomics*, **18**, 1-18.

- 1252 Wang J, Street NR, Scofield DG, Ingvarsson PK (2016) Natural selection and recombination
1253 rate variation shape nucleotide polymorphism across the genomes of three related *Populus*
1254 species. *Genetics*, **202**, 1185-1200.
- 1255 Warr A, Robert C, Hume D, *et al.* (2015) Exome sequencing : Current and Future
1256 perspectives. *Genes, Genomes, Genetics*, **5**, 1543-1550.
- 1257 Weir BS, Cockerham CC (1984) Estimating *F*-statistics for the analysis of population
1258 structure. *Evolution*, **38**, 1358–1370.
- 1259 Weir BS, Hill WG (2002) Estimating *F*-statistics. *Annual Review of Genetics*, **36**, 721–750.
- 1260 Willing E-M, Dreyer C, van Oosterhout C (2012) Estimates of genetic differentiation
1261 measured by F_{ST} do not necessarily require large sample sizes when using many SNP
1262 markers. PLoS ONE, **7**, e42649.
- 1263 Wolf JB, Ellegren H (2017) Making sense of genomic islands of differentiation in light of
1264 speciation. *Nature Reviews Genetics*, **18**, 87–100.
- 1265 Wu CI (2001) The genic view of the process of speciation. *Journal of Evolutionary Biology*
1266 **14**, 851–865.
- 1267 Zanetto A, Roussel G, Kremer A (1994) Geographic variation of inter- specific differentiation
1268 between *Quercus robur* L. and *Quercus petraea* (Matt.) Liebl. *Forest Genetics*, **1**, 111-123.
- 1269 Zhang M, Zhou L, Bawa R, Suren H, Holliday JA (2016) Recombination rate variation,
1270 hitchhiking, and demographic history shape deleterious load in poplar. *Molecular Biology and*
1271 *Evolution*, **33**, 2899–2910.
- 1272 Zhou L, Bawa R, Holliday JA. 2014. Exome resequencing reveals signatures of demographic
1273 and adaptive processes across the genome and range of black cottonwood (*Populus*
1274 *trichocarpa*). *Molecular Ecology*, **23**, 2486–2499.
- 1275 Zhou L, Holliday JA (2012) Targeted enrichment of the black cottonwood (*Populus*
1276 *trichocarpa*) gene space using sequence capture. *BMC Genomics*, **13**, 703.

1277

1278 **Author contributions**

1279 Funding acquisition: AK, PGG, CP, and MLDL; Initial conception and individuals sampling:
1280 PGG, AK, CP, MPR, VL; Bioinformatics strategy and experimental design: PGG, TL; DNA
1281 extraction and quality check: VL; Sequence Data acquisition: PGG, CP, TL, VL; Individuals'
1282 identification checks for quality control VL, CL, PL; Pilot study: VL, PGG; Working
1283 assembly: JMF, PGG; Primer design and amplicon choice: PGG, VL, TD; Original candidate
1284 gene lists choice: PGG, TL, JMF, CP, AK, TD, CR, MLDL, GLP, ChB, EG, CaB, NT, PA;

1285 Bioinformatics tools: TL and PGG (SeqQual pipeline and R scripts), JMF and AF (Bioperl
1286 and R scripts), PA, CL, VelM, JT, FH, TD (SeqQual tests), FR (website); Visual
1287 Chromatogram checks, SNP/assembly validations: PGG, VL, TD, PA, TL, MLDL, CaB,
1288 ChB, CL, CR and EG; Bioinformatic and population genetic analyses: PGG, TL, SM, ChB:
1289 Functional annotation: TL, PGG, VelM, PA; Manuscript draft: PGG; Manuscript review and
1290 edition: PGG, SM, CL, ChB, TL; all authors agreed on the manuscript.
1291

1292 **Supporting Information**

1293 **Fig. S1** Sampling site locations within the natural geographic distribution of *Q. petraea* and
1294 *Q. robur*. Vector map is from <http://www.naturalearthdata.com> and distribution areas from
1295 Euforgen (<http://www.euforgen.org/distribution-maps/>)

1296 **Fig. S2** Working assembly steps and softwares (A), and bioinformatic strategy for search of
1297 candidate genes and amplicon choice (B).

1298 **Fig. S3** Plots of the ΔK values from the Evanno *et al.* (2005) method (S3-A, -B, -C, -D, -E),
1299 and of the mean values of the estimated probability \ln (of the data given K) with standard
1300 deviations for K ranging from 1 to 5 (S3-F to S3-J), which show support for $K=2$. Plots are
1301 from the STRUCTURE HARVESTER program.

1302 **Fig. S4** Distributions of Gene Ontology (GO) terms for the consensus sequences in Appendix
1303 S3, at level 2 (-A, -B, -C) and level 3 (-D, -E, -F): A- and D- for Biological Process, B- and E-
1304 for Molecular Function, C- and F- for Cellular Component. Annotation rules: E-value $<10^{-30}$,
1305 annotation cut-off 70, GO weight 5, HSP coverage cutoff 33%. Filtering applies for at least 5
1306 sequences and a node score of 5 per GO term (but see rare exceptions in Table S2).

1307 **Fig. S5** Distributions of GO terms across different gene lists (*bud*, *abiotic* and *biotic*) at
1308 Biological Level 2, and Fisher exact tests across pairs of sequence clusters with the same GO
1309 terms between the random list and other lists. Significance levels *: $P<0.05$.

1310 **Fig. S6-A to S6-J** Posterior assignment probabilities (Q -values) of 24 individuals attributed to
1311 2 clusters (STRUCTURE analysis) for different numbers of polymorphisms, different sampling
1312 of SNP data, and different plots of credible intervals.

1313 **Fig. S7** Mean number of new variants brought by each new distinct individual added to all
1314 possible initial sample size combinations (-A); Number of high-quality variant positions per
1315 100 base pair (bp) across 852 gene fragments ranked by their length (bp), overall and for each
1316 species (-B).

1317 **Fig. S8** Comparison of nucleotide diversity (θ .pi) distributions between main cpDNA
1318 lineages (B and A or C) for *Q. robur* (586 genes) and *Q. petraea* (449 genes). The histogram
1319 represents lineage B for *Q. robur*. Data are available in both lineages within each species for
1320 at least 8 gametes per lineage, and a minimum of 200 bp per gene fragment.

1321 **Table S1** Description of amplicons: primer sequences, original candidate gene list, targeted
1322 biological functions (see references), candidate gene type, fragment expected size and
1323 position in the *orict* original working assembly, preliminary results based nucleotide quality
1324 for obtained sequences, and validation decision after excluding paralog amplification.

1325 **Table S2** Functional annotation results from Blast2GO (-A), comparison of BlastX best hits
1326 results (according to E -values) between consensus sequences of the *orict* working assembly
1327 and the *ocv4* assembly (-B), and comparisons of BlastN results of consensus sequence for
1328 both *orict* and corresponding expected amplicon (*orict-cut*) onto *ocv4* (-C).

1329 **Table S3** Description of all variants single base positions, with sample sizes, alleles,
1330 genotypes counts, various statistics, and generic format for genotyping essays input data.
1331 Species samples exclude the 2 most introgressed individuals.

1332 **Table S4** Description of all polymorphisms as in Table S3, but with a characterization of the
1333 length, sequence motifs, contiguous base positions for complex polymorphic regions
1334 including indels, SNPs and SSRs (see also Table S5 for SSR positions).

- 1335 **Table S5** SSR patterns as detected from the *mreps* software.
- 1336 **Appendix S1** Additional method details.
- 1337 **Appendix S2** Contigs of the original working assembly used for selecting candidate gene
1338 regions and design amplicon primers, including consensus sequences and reads where
1339 nucleotides with Phred score below 20 have been masked.
- 1340 **Appendix S3** Sequences of chosen contig consensus and singletons sequences for functional
1341 annotation analyses.
- 1342 **Appendix S4** Consensus sequences of 852 genomic regions obtained in this study for
1343 *Quercus petraea* and *Q. Robur* individuals. “(N)⁹” : represents a low-quality fragment of a
1344 length below ~1 kb separating Forward and Reverse amplicons; “n” represents positions with
1345 a majority of nucleotides with phd score below 30. “(-)^x”: means that the insertion is a minor
1346 allele at that position, x being the size of the indel.
- 1347 **Appendix S5** Nucleotide sequence data of 394 gene regions for one *Quercus ilex* individual,
1348 heterozygote sites being indicated by IUPAC codes.
- 1349 **Appendix S6** Outputs from Blast2GO analyses.