

# 1 **Low Rate of Somatic Mutations in a Long-Lived Oak Tree**

2 Namrata Sarkar<sup>1,2,3,9</sup>, Emanuel Schmid-Siegert<sup>4,9</sup>, Christian Iseli<sup>4,9</sup>, Sandra Calderon<sup>4</sup>,  
3 Caroline Gouhier-Darimont<sup>5</sup>, Jacqueline Chrast<sup>1</sup>, Pietro Cattaneo<sup>5</sup>, Frédéric Schütz<sup>1</sup>, Laurent  
4 Farinelli<sup>6</sup>, Marco Pagni<sup>4</sup>, Michel Schneider<sup>7</sup>, Jérémie Voumard<sup>8</sup>, Michel Jaboyedoff<sup>8</sup>,  
5 Christian Fankhauser<sup>1\*</sup>, Christian S. Hardtke<sup>5\*</sup>, Laurent Keller<sup>2\*</sup>, John R. Pannell<sup>2\*</sup>,  
6 Alexandre Reymond<sup>1\*</sup>, Marc Robinson-Rechavi<sup>2,3\*</sup>, Ioannis Xenarios<sup>1,4,7\*</sup>, Philippe  
7 Reymond<sup>5,10\*</sup>

8

9 <sup>1</sup>Center for Integrative Genomics, University of Lausanne, 1015 Lausanne, Switzerland

10 <sup>2</sup>Department of Ecology and Evolution, University of Lausanne, 1015 Lausanne, Switzerland

11 <sup>3</sup>Evolutionary Bioinformatics Group, Swiss Institute of Bioinformatics, 1015 Lausanne,  
12 Switzerland

13 <sup>4</sup>Vital-IT Competence Center, Swiss Institute of Bioinformatics, 1015 Lausanne, Switzerland

14 <sup>5</sup>Department of Plant Molecular Biology, University of Lausanne, 1015 Lausanne,  
15 Switzerland

16 <sup>6</sup>Fasteris SA, 1228 Plan-les-Ouates, Switzerland

17 <sup>7</sup>Swiss-Prot group, Swiss Institute of Bioinformatics, 1211 Geneva, Switzerland

18 <sup>8</sup>Risk Analysis Group, Institute of Earth Sciences, University of Lausanne, 1015 Lausanne,  
19 Switzerland

20 <sup>9</sup>These authors contributed equally

21 <sup>10</sup>Lead contact

22 \*Correspondence: christian.fankhauser@unil.ch, christian.hardtke@unil.ch,

23 laurent.keller@unil.ch, john.pannell@unil.ch, alexandre.reymond@unil.ch, marc.robinson-

24 rechavi@unil.ch, ioannis.xenarios@unil.ch, philippe.reymond@unil.ch

25

26

27

## 28 **SUMMARY**

29 Because plants do not possess a proper germline, deleterious mutations that occur in the soma  
30 can be passed to gametes. It has generally been assumed that the large number of somatic cell  
31 divisions separating zygote from gamete formation in long-lived plants should lead to many  
32 mutations. However, a recent study showed that surprisingly few cell divisions separate  
33 apical stem cells from axillary stem cells in annual plants, challenging this view. To test this  
34 prediction, we generated and analysed the full genome sequence of two terminal branches of  
35 a 234-year-old oak tree and found very few fixed somatic single-nucleotide variants (SNVs),  
36 whose sequential appearance in the tree could reliably be traced back along nested sectors of  
37 younger branches. Our data indicate that the stem cells of shoot meristems in trees are  
38 robustly protected from accumulation of mutations, analogous to the germline in animals.

39

40

## 41 **INTRODUCTION**

42 Accumulation of deleterious mutations is a fundamental parameter in plant ageing and  
43 evolution [1, 2]. Because the pedigree of cell division that generates somatic tissue is poorly  
44 understood, the number of cell divisions that separate zygote from gamete formation is  
45 difficult to estimate; this number is expected to be particularly large in trees and could in  
46 theory lead to a large number of DNA replication errors [3-5]. Tree architecture is determined  
47 by the modular growth of apical meristems, which contain stem cells. These cells divide and  
48 produce progenitor cells that undergo division, elongation and differentiation to form a  
49 vegetative shoot, the branch. Axillary meristems are formed at the base of leaf axils and are

50 responsible for the emergence of side branches. They are separated from apical meristems by  
51 elongating internodes. In oak, early and repeated growth cessation of terminal apical  
52 meristems leads to a branching pattern originating from such axillary meristems. In turn,  
53 axillary meristems grow out and produce secondary axillary meristems. This process is  
54 reiterated indeterminately to produce highly ramified trees of large stature, resulting in  
55 thousands of terminal ramets [6, 7].

56 Classical studies of shoot apical meristem organization have reported that the most  
57 distal zone has a significantly lower rate of cell division than more basal regions of the apex,  
58 and might therefore be relatively protected from replication errors [8, 9]. In a recent study  
59 that followed the fate of dividing cells in the apical meristems of *Arabidopsis thaliana* and  
60 tomato, Burian et al. [10] showed that an unexpectedly low number of divisions separate  
61 apical from axillary meristems. In these herbaceous plants, axillary meristems are separated  
62 from apical meristem stem cells by seven to nine cell divisions, with internode growth  
63 occurring through the division of cells behind the meristem. The number of cell divisions  
64 between early embryonic stem cells and terminal meristems thus depends more on the  
65 number of branching events than on absolute plant size. Moreover, only three or four semi-  
66 permanent apical stem cells give rise to independent sectors of the growing shoot, so that a  
67 mutation in one apical stem cell may propagate and be fixed in a domain of the apical shoot  
68 that contains an axillary meristem [10]. Alternatively, mutations that arise in sub-apical stem  
69 cells may be partly fixed in axillary meristems and are either lost or fully fixed in secondary-  
70 order axillary meristems.

71 In trees, because axillary meristems replace apical meristems iteratively, the  
72 cumulative number of cell divisions separating meristems determines the rate of genetic  
73 aging and the potential accumulation of somatic mutations. If the same growth pattern  
74 described above for *A. thaliana* and tomato applies to trees, their somatic mutation rate might

75 be much lower than is commonly thought, and the majority of fixed mutations should be  
76 found in relatively small sectors as nested sets of mutations. To test these predictions, we  
77 sequenced the full genome of two terminal branches of an iconic old oak tree (*Quercus*  
78 *robur*).

79

## 80 **RESULTS AND DISCUSSION**

### 81 **Napoleon Oak Genome Sequencing**

82 We conducted our study on an oak tree known as the ‘Napoleon Oak’ by the academic  
83 community of the University of Lausanne. The tree was 22 years old when, on May 12, 1800,  
84 Napoleon Bonaparte and his troops crossed what is now the Lausanne University campus, on  
85 their way to conquer Italy. It was standing next to the road, or may have been transplanted in  
86 honour of the future emperor's passage, whence its campus nickname. At the time of sample  
87 harvest for our study, the dividing apical meristems of this magnificent tree (Figure 1, Figure  
88 S1) had been exposed for 234 years to potential environmental mutagens, such as UV and  
89 radioactive radiation.

90 To identify fixed somatic variants (i.e., those present in an entire sector of the  
91 Napoleon Oak) and to reconstruct their origin and distribution among branches, we collected  
92 26 leaf samples from different locations on the tree. We first sequenced the genome from  
93 leaves sampled on terminal ramets of one lower and one upper branch of the tree. We then  
94 used a combination of short-read Illumina and single-molecule real-time (SMRT, Pacific  
95 Biosciences) sequencing to generate a *de novo* assembly of the oak genome. After removing  
96 contigs <1000bp, we established a draft sequence of ca. 720 megabases (Mb) at a coverage of  
97 ca. 70X, with 85,557 scaffolds and a N50 length of 17,014. Our sequence is thus in broad  
98 agreement with the published estimated genome size of 740 Mbp [11]. The oak genome  
99 encodes 49,444 predicted protein-coding loci (Table 1).

## 100 **Identification of Somatic Mutations Between Genomes from Two Branches**

101 We used two approaches to identify SNVs (single-nucleotide variants) between the  
102 sequenced genomes of the two terminal branches of our initial sample, for regions of both  
103 high and low coverage. In total, reads from ca. 650 Mb of non-Ns containing sequence could  
104 be assessed. First, we aligned Illumina paired-reads on the repeat-masked genome in  
105 combination with the GATK [12] variant caller. This allowed us to establish a list of 314,865  
106 potential SNV candidates. On the basis of a confidence score  $\geq 300$  on the heterozygous sites  
107 and  $\geq 200$  on homozygous sites, we further selected 1,536 putative SNVs for experimental  
108 validation by both PCR-seq [13] and Sanger sequencing. Of the 1,536 candidates, only seven  
109 could be confirmed (see Experimental Procedures).

110         Second, we used fetchGWI [14] to map read pairs to the entire non-masked genome.  
111 We were able to call 5,330 potential SNVs from the mapped reads using a simple read pileup  
112 process followed by detection of positions where the pileup differed from the reference  
113 genome. We systematically browsed candidate positions to evaluate the quality of the  
114 mapping in the surrounding region, and to discriminate between well-assembled high-quality  
115 regions with two alleles per sample and potentially poorly assembled repeated regions. This  
116 approach allowed us to select 82 putatively variable positions, including the seven already  
117 identified using the repeat-masked genome analysis described above. Ten of the remaining 75  
118 candidates identified using the second approach were confirmed by PCR-seq and Sanger  
119 sequencing, increasing the total number of confirmed SNVs separating the two genomes to  
120 17 (Figure 1, Table S1). Based on a conservative estimate, we are likely to have missed no  
121 more than 17 further such sites (see Experimental Procedures). All confirmed SNVs were  
122 heterozygous, as expected for novel somatic mutations. Intriguingly, two SNVs were found  
123 on the same contig, separated by only 12 bp (Figure 1, Table S1).

124

## 125 **Nested Distribution of SNVs Throughout the Napoleon Oak**

126 Having confidently established 17 SNVs, we then assessed their occurrence throughout the  
127 tree. We used Sanger sequencing to genotype the remaining 24 terminal branches sampled  
128 from other parts of the tree and checked for the presence of each SNV. As might be expected,  
129 the SNVs were found in different sectors of the tree in a nested hierarchy that clearly  
130 indicates the accumulation of mutations along branches during development (Figure 1, Figure  
131 S2). Specifically, SNV1 and SNV2 were located in a large sector of connected branches that  
132 are distributed in the upper half of the tree (Figure 1). SNV3 was restricted to a portion of the  
133 same sector, whereas SNV4 to SNV10 were found at the top of one of the two branches for  
134 which the full genome had been sequenced. Similarly, SNV11 to SNV14 were restricted to  
135 the terminal fork of a principle lower branch from which the other genome had been fully  
136 sequenced, whereas SNV15 to SNV17 were present in only one sector arising from this fork  
137 (Figure 1, Figure S2). These results both provide independent confirmation of the originally  
138 identified SNVs, and demonstrate their gradual, nested appearance and fixation in  
139 developmentally connected branches during growth of the oak tree. Thus, while the exact  
140 ontogeny of the Napoleon Oak may be difficult to reconstruct, our SNV analysis generated a  
141 nested set of lineages supported by derived mutations, analogous to a phylogenetic tree.

142

## 143 **Somatic Mutation Rate is Low in the Napoleon Oak**

144 The spontaneous mutation rate in plants has been estimated to range from  $5 \times 10^{-9}$  to  $30 \times 10^{-9}$   
145 substitutions/site/generation, based on mutations accumulated during divergence between  
146 monocots and dicots [15] or divergence between independent lines of *A. thaliana* maintained  
147 in the laboratory [16, 17]. *A. thaliana* is an annual plant that reaches approximately 30 cm in  
148 height before producing seeds. In contrast, the physical distance traced along branches  
149 between the terminal branches we sequenced for the Napoleon Oak is about 40 m (Figure 1).

150 Assuming similar cell sizes between oak and *A. thaliana*, there should have been about 133  
151 (4,000 cm/ 30 cm) times more mitotic divisions separating the extremes of somatic lineages  
152 in oaks than in *A. thaliana*. Under the assumption that the per-generation mutation rate is  
153 correlated with the number of mitotic divisions from zygote to gametes of the next generation  
154 [3, 8], the mitotic mutation rate for the Napoleon Oak should be between  $6.6 \times 10^{-7}$  and  $4.0 \times$   
155  $10^{-6}$  substitution/site/generation. With such a somatic mutation rate, there should be between  
156 433 and 2,600 SNVs across the 650 Mb sequenced between branch ends of the tree, or  
157 conservatively between 305 and 1,832 SNVs across the 458,143,725 nucleotides with a read  
158 coverage  $\geq 8$  in both samples. These values are much higher than the 17 SNVs we actually  
159 found. Even under the conservative assumption that we missed 17 other SNVs from the list  
160 of potential variants and between 12.9 and 79 false negatives (see Experimental Procedures),  
161 a higher-end estimate of the total number of oak SNVs ranges from 46.9 to 113. Both values  
162 are still substantially (and significantly) lower than the expected values (305 to 1,832)  
163 (Fisher's exact test, both  $P < 2.2E-16$ ). Our finding seems to imply either much lower  
164 generational or mitotic mutation rates than generally assumed, or a different model for the  
165 accumulation of somatic mutations in trees, such as that proposed by [10]. Although the oak  
166 lineages sampled have not been separated by any meiosis events, which in yeast was found to  
167 elevate the generational mutation rate [19], they have been exposed to the natural  
168 environment, which in *A. thaliana* is known to significantly enhance mutation rate when  
169 compared to a controlled lab environment [20].

170 Most of SNVs identified in the Napoleon Oak were found only in single branches.  
171 Early sequestration of axillary meristems has been found in *A. thaliana* and tomato, two  
172 herbaceous plants that differ markedly in the developmental fate of their apical meristems  
173 [10]. Whereas *A. thaliana* forms a rosette of leaves and extends an inflorescence shoot at the  
174 time of flowering, the tomato apical meristems are separated by elongated internodes. Oak

175 trees have apical meristems of similar diameters to those of tomato [10, 21, and Figure 2] and  
176 show similar ontogeny. It thus seems reasonable to suppose that the growth pattern described  
177 in *A. thaliana* and tomato is quite general in flowering plants and might also apply to long-  
178 lived trees. If so, it would be consistent with the low number of SNVs we have identified in  
179 the oak.

180 We found that G:C→A:T transitions were the most frequent class of SNVs observed  
181 in the oak (Figure 3). Ultraviolet (UV) light causes G:C→A:T transitions at dipyrimidine  
182 sites in plants [22]. Among the 11 G:C→A:T transitions that we observed, seven were in a  
183 dipyrimidine context (Table S1). In addition, spontaneous deamination of methylated  
184 cytosine leads to thymine change at CpG or CpNG sites [22]. However, there were only three  
185 G:C→A:T transitions in such a context (Table S1). It thus seems plausible that UV light may  
186 have caused most of the G:C→A:T transitions we observed, although other factors, such as  
187 cytosine deamination and replication errors, may account for other SNVs.

188 Taken together, our results suggest that mutations due to replication errors in trees  
189 may be less important than environmentally induced mutations, as hypothesized by Burian et  
190 al. [10]. Mutations accumulate with age, irrespective of plant stature, and long-term exposure  
191 to UV radiation contributes to such changes. However, oaks protect their meristems in buds  
192 under multi-layered leaf-like structures (Figure 2), potentially reducing the incidence of UV  
193 mutagenesis. If so, the surprisingly low somatic mutation rate implied by our data would be  
194 consistent both with the pattern of cell division hypothesised by Burian et al. [10], as well as  
195 with the protective nature of oak bud morphology. In this context, it is noteworthy that there  
196 was no evidence for an expansion of DNA-repair genes in the oak genome compared to *A.*  
197 *thaliana* (Table S2).

198 Sixteen of the 17 SNVs identified in the Napoleon Oak occurred in introns or non-  
199 coding sequences that are probably neutral. The remaining one (SNV1), which occurred in a



200 large sector of the tree, generates an arginine-to-glycine conversion in a putative E3-ubiquitin  
201 ligase (Table S1). The functional impact of exchanging a positively charged arginine with a  
202 non-charged and smaller glycine residue is unknown and deserves further analysis.

203 To our knowledge, only two examples of functional mosaicism have been reported in  
204 trees, a low incidence that might be attributable to the low somatic mutation rate implied by  
205 our analysis. Although most non-neutral mutations should be maladaptive, eucalyptus trees  
206 have been observed with a few branches that are biochemically distinct from the rest of the  
207 canopy and have become resistant to Christmas beetle defoliation [23, 24]. Functionally  
208 relevant somatic mutations, such as SNV1 in our study, may thus occasionally contribute to  
209 adaptive evolution if transferred to the fruits, but will more typically increase the genetic load  
210 of a population, with implications for inbreeding depression and mating-system evolution.

211 Our results throw new light on explanations proposed for differences in the  
212 distribution of mating systems between short- and long-lived plants. While many annuals and  
213 short-lived plants have undergone evolutionary transitions from outcrossing to selfing [25],  
214 often involving a loss of self-incompatibility systems [26], long-lived woody species are  
215 more likely to be fully outcrossing [27], including oaks [28]. Theoretical analysis indicates  
216 that a high somatic mutation rate could account for this difference, because somatic  
217 mutations would contribute to the genetic load of the population and thus to inbreeding  
218 depression, disfavouring self-fertilization [3]. Inbreeding depression is indeed higher in long-  
219 lived woody species than annuals [29], and the observation of higher inbreeding depression  
220 caused by within-branch than between-branch selfing points to the accumulation of different  
221 deleterious somatic mutations in different sectors of the plant [5]. However, our finding now  
222 challenges the notion that the breeding system of long-lived trees is constrained by a high rate  
223 of somatic mutations.

224 The results of our study, in conjunction with those of Burian et al. [10], have  
225 important implications for how we should view one of the most fundamental ways in which  
226 plants differ from animals – their absence of a germline. In oak, iterative growth of axillary  
227 meristems produces terminal branches that carry stem cells. As in other plants, favourable  
228 conditions induce stem cells to produce floral buds and ultimately the gametes of the next  
229 generation. These stem cells are functionally analogous to germ cells in metazoans and result  
230 from a limited number of divisions that prevent an accumulation of replicative errors.

231 Our data give an unprecedented view on the limited role played by somatic mutations  
232 in long-lived plants, and support the view that stem cells in trees, although vulnerable to  
233 environment-induced mutations, are probably quite well protected from them. Consistent  
234 with this finding, a recent study in *A. thaliana* has shown that the number of cell divisions  
235 from germination to gametogenesis is independent of life span and vegetative growth [30].  
236 Our study also illustrates the potential for analyses of multiple genomes from single  
237 individuals, which throw exciting new light on the rate, distribution and potential impact of  
238 somatic mutations in both plant and animal tissues [31-34].

239

## 240 **SUPPLEMENTAL INFORMATION**

241 Supplemental Information includes five figures and seven tables.

242

## 243 **AUTHOR CONTRUBUTIONS**

244 L. F. sequenced the genome. E.S.-S., S.C., M.P. assembled and annotated the genome. N.S.,  
245 E.S.-S., C.I. identified SNVs. C.G.-D., J.C. extracted DNA and confirmed SNVs. E.S.-S.,  
246 M.R.-R. analyzed genome duplication. P.C. produced cross-sections of oak apical meristems.  
247 M. S. established a list of DNA repair genes. F. S. provided statistical help with the analyses.

248 J.V., M.J. produced a 3D model of the oak tree. C.H., C.F., L.K., I.X., M.R.-R., J.P., A.R.,  
249 P.R. conceived the project and wrote the manuscript.

250

## 251 **ACKNOWLEDGMENTS**

252 This work was funded by the University of Lausanne through a supportive grant from the  
253 University rectorate and by the Swiss National Science Foundation (Agora Grant  
254 CRAGI3\_145652). The Pacific Biosciences RS II sequencing was performed at the Lausanne  
255 Genomic Technologies Facility (GTF). The purchase of the GTF's RS II instrument was  
256 financed in part by the *Loterie Romande* through the *Fondation pour la Recherche en*  
257 *Médecine Génétique*. We thank Keith Harshman, Johann Weber, and Mélanie Dupasquier  
258 from the GTF for sequencing. We thank Cris Kuhlemeier for sharing unpublished results,  
259 Jean Tercier for tree-ring analysis, Woodtli+Leuba SA for sample collection, and Jean-  
260 Jacques Strahm and Marco Bonetti for providing oak images. All Illumina reads and SMRT  
261 sequences have been deposited in GenBank under accession BioProject PRJNA327502.

262

## 263 **REFERENCES**

- 264 1. Klekowski, E.J., and Godfrey, P.J. (1989). Ageing and mutation in plants. *Nature* 340,  
265 389-391.
- 266 2. Schultz, S.T., Lynch, M., and Willis JH. (1999). Spontaneous deleterious mutation in  
267 *Arabidopsis thaliana*. *Proc. Natl. Acad. Sci. USA* 96, 11393–11398.
- 268 3. Scofield, D.G., and Schultz, S.T. (2006). Mitosis, stature and evolution of plant mating  
269 systems: low- $\Phi$  and high- $\Phi$  plants. *Proc. Royal Soc. London B: Biol. Sci.* 273:275–  
270 282.
- 271 4. Ally., D., Ritland, K., and Otto, S.P. (2010). Aging in a long-lived clonal tree. *PLoS*  
272 *Biol.* 8, e1000454.
- 273 5. Bobiwash, K., Schultz, S.T., and Schoen, D.J. (2013). Somatic deleterious mutation  
274 rate in a woody plant: estimation from phenotypic data. *Heredity* 111, 338–344.
- 275 6. Halle, F. (1986). Modular growth in seed plants. *Phil. Trans. Royal Soc. London B:*  
276 *Biol. Sci.* 313, 77–87.
- 277 7. Millet, J. (2012). L'architecture des arbres des régions tempérées: son histoire, ses  
278 concepts, ses usages. (MultiMondes).
- 279 8. Romberger, J.A., Hejnowicz, Z., and Hill, J.F. (1993). Plant structure: function and  
280 development. (Springer-Verlag).

- 281 9. Kwiatkowska, D. (2008). Flowering and apical meristem growth dynamics. *J. Exp. Bot.*  
282 *59*, 187-201.
- 283 10. Burian, A., Barbier de Reuille, P., and Kuhlemeier, C. (2016). Patterns of stem cell  
284 divisions contribute to plant longevity. *Curr. Biol.* *26*, 1385-1394.
- 285 11. Plomion, C., Aury, J.M., Amsellem, J., Alaeitabar, T., Barbe, V., Belser, C., Bergès, H.,  
286 Bodénès, C., Boudet, N., Boury, C., et al. (2016). Decoding the oak genome: public  
287 release of sequence data, assembly, annotation and publication strategies. *Mol. Ecol.*  
288 *Resour.* *16*, 254-265.
- 289 12. McKenna, A., Hanna, M., Banks, E., Sivachenko, A., Cibulskis, K., Kernytzky, A.,  
290 Garimella, K., Altshuler, D., Gabriel, S., Daly, M., et al. (2010). The genome analysis  
291 toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data.  
292 *Genome Res.* *20*, 1297-303.
- 293 13. Howald, C., Tanzer, A., Chrast, J., Kokocinski, F., Derrien, T., Walters, N., Gonzalez,  
294 J.M., Frankish, A., Aken, B.L., Hourlier, T., et al. (2012). Combining RT-PCR-seq to  
295 catalog all genic elements encoded in the human genome. *Genome Res.* *22*, 1698-1710.
- 296 14. Iseli, C., Ambrosini, G., Bucher, P., and Jongeneel, C.V. (2007). Indexing strategies for  
297 rapid searches of short words in genome sequences. *PLoS One* *2*, e579.
- 298 15. Wolfe, K.H., Li, W.H., and Sharp, P.M. (1987). Rates of nucleotide substitution vary  
299 greatly among plant mitochondrial, chloroplast, and nuclear DNAs. *Proc. Natl. Acad.*  
300 *Sci. USA* *84*, 9054–9058.
- 301 16. Ossowski, S., Schneeberger, K., Lucas-Lledó, J.I., Warthmann, N., Clark, R.M., Shaw,  
302 R.G., Weigel, D., and Lynch, M. (2010). The rate and molecular spectrum of  
303 spontaneous mutations in *Arabidopsis thaliana*. *Science* *327*, 92–94.
- 304 17. Yang, S., Wang, L., Huang, J., Zhang, X., Yuan, Y., Chen, J.Q., Hurst, L.D., and Tian,  
305 D. (2015). Parent–progeny sequencing indicates higher mutation rates in heterozygotes.  
306 *Nature* *523*, 463-467.
- 307 18. Scofield, D.G. (2006). Medial pith cells per meter in twigs as a proxy for mitotic  
308 growth rate ( $\Phi/m$ ) in the apical meristem. *Am. J. Bot.* *93*, 1740-1747.
- 309 19. Rattray, A., Santoyo, G., Shafer, B., and Strathern, J.N. (2015). Elevated mutation rate  
310 during meiosis in *Saccharomyces cerevisiae*. *PLoS Genet.* *11*, e1004910.
- 311 20. Rutter, M.T., Shaw, F.H., and Fenster, C.B. (2010). Spontaneous mutation parameters  
312 for *Arabidopsis thaliana* measured in the wild. *Evolution* *64*, 1825-1835.
- 313 21. Barton, M.K. (2010). Twenty years on: The inner workings of the shoot apical  
314 meristem, a developmental dynamo. *Dev. Biol.* *341*, 95-113.
- 315 22. Britt, A.B. (1996). DNA damage and repair in plants. *Annu. Rev. Plant Physiol. Plant*  
316 *Mol. Biol.* *47*, 75–100.
- 317 23. Edwards, P.B., Wanjura, W.J., Brown, W.V., and Dearn, J.M. (1990). Mosaic  
318 resistance in plants. *Nature* *347*, 434.
- 319 24. Padovan, A., Lanfear, R., Keszei, A., Foley, W.J., and Külheim, C. (2013). Differences  
320 in gene expression within a striking phenotypic mosaic *Eucalyptus* tree that varies in  
321 susceptibility to herbivory. *BMC Plant Biol.* *13*, 29.
- 322 25. Stebbins, G.L. (1950). *Variation and evolution in plants*. (Columbia University Press).
- 323 26. Goldberg, E.E., Kohn, J.R., Lande, R., Robertson, K.A., Smith, S.A., and Igić, B.  
324 (2010). Species selection maintains self-incompatibility. *Science* *330*, 493-495.
- 325 27. Barrett, S.C.H., Harder, L.D., and Worley, A.C. (1996). The comparative biology of  
326 pollination and mating in flowering plants. *Phil. Tran. Royal Soc. London B: Biol. Sci.*  
327 *351*, 1271-1280.
- 328 28. Streiff, R., Ducouso, A., Lexer, C., Steinkellner, H., Gloessl, J. and Kremer, A. (2009).  
329 Pollen dispersal inferred from paternity analysis in a mixed oak stand of *Quercus robur*  
330 *L.* and *Q. petraea* (Matt.) Liebl. *Mol. Ecol.* *8*, 831-841.

- 331 29. Goodwillie, C., Kalisz, S., and Eckert, C.E. (2005). The evolutionary enigma of mixed  
332 mating systems in plants: occurrence, theoretical explanations, and empirical evidence.  
333 *Annu. Rev. Ecol. Evo. Syst.* *36*, 47-79.
- 334 30. Watson, J.M., Platzer, A., Kazda, A., Akimcheva, S., Valuchova, S., Nizhynska, V.,  
335 Nordborg, M., and Riha, K. (2016). Germline replications and somatic mutation  
336 accumulation are independent of vegetative life span in *Arabidopsis*. *Proc. Natl. Acad.*  
337 *Sci. USA* *113*, 12226–12231.
- 338 31. James, R.L. (2013). Genome mosaicism - One human, multiple genomes. *Science* *341*,  
339 358–359.
- 340 32. Behjati, S., Huch, M., van Boxtel, R., Karthaus, W., Wedge, D.C., Tamuri, A.U.,  
341 Martincorena, I., Petljak, M., Alexandrov, L.B., Gundem, G., et al. (2014). Genome  
342 sequencing of normal cells reveals developmental lineages and mutational processes.  
343 *Nature* *513*, 422–425.
- 344 33. Lodato, M.A., Woodworth, M.B., Lee, S., Evrony, G.D., Mehta, B.K., Karger, A., Lee,  
345 S., Chittenden, T.W., D'Gama, A.M., Cai, X., et al. (2015). Somatic mutation in single  
346 human neurons tracks developmental and transcriptional history. *Science* *350*, 94–98.
- 347 34. Martincorena, I., Roshan, A., Gerstung, M., Ellis, P., Van Loo, P., McLaren, S.,  
348 Wedge, D.C., Fullam, A., Alexandrov, L.B., Tubio, J.M., et al. (2015). Tumor  
349 evolution. High burden and pervasive positive selection of somatic mutations in normal  
350 human skin. *Science* *348*, 880–886.
- 351 35. Bolger, A.M., Lohse, M., and Usadel, B. (2014). Trimmomatic: A flexible trimmer for  
352 Illumina Sequence Data. *Bioinformatics* *30*, 2114-2120.
- 353 36. Luo, R., Liu, B., Xie, Y., Li, Z., Huang, W., Yuan, J., He, G., Chen, Y., Pan, Q., Liu,  
354 Y., et al. (2012). SOAPdenovo2: an empirically improved memory-efficient short-read  
355 de novo assembler. *GigaScience* *1*, 18.
- 356 37. Boetzer, M., Henkel, C.V., Jansen, H.J., Butler, D., and Pirovano, W. (2011).  
357 Scaffolding pre-assembled contigs using SSPACE. *Bioinformatics* *27*, 578–579.
- 358 38. Boetzer, M., and Pirovano, W. (2012). Toward almost closed genomes with GapFiller.  
359 *Genome Biol.* *13*, R56.
- 360 39. Segata, N., Waldron, L., Ballarini, A., Narasimhan, V., Jousson, O., and Huttenhower,  
361 C. (2012). Metagenomic microbial community profiling using unique clade-specific  
362 marker genes. *Nature Methods* *9*, 811-814.
- 363 40. Altschul, S.F., Gish, W., Miller, W., Myers, E.W., and Lipman, D.J. (1990). Basic local  
364 alignment search tool. *J. Mol. Biol.* *215*, 403-410.
- 365 41. English, A.C., Richards, S., Han, Y., Wang, M., Vee, V., Qu, J., Qin, X., Muzny, D.M.,  
366 Reid, J.G., Worley, K.C., et al. (2012). Mind the gap: Upgrading genomes with Pacific  
367 Biosciences RS Long-Read Sequencing Technology. *PLoS One* *7*, 11.
- 368 42. Stanke, M., Steinkamp, R., Waack, S., and Morgenstern, B. (2004). AUGUSTUS: a  
369 web server for gene finding in eukaryotes. *Nucl. Acids Res.* *32*, W309-312.
- 370 43. Tran, V.D., De Coi, N., Feuermann, M., Schmid-Siegert, E., Bağcı, E.T., Mignon, B.,  
371 Waridel, P., Peter, C., Pradervand, S., Pagni, M., et al. (2016). RNA sequencing-based  
372 genome reannotation of the dermatophyte *Arthroderma benhamiae* and characterization  
373 of its secretome and whole gene expression profile during infection. *mSystems* *11*,  
374 e00036-16.
- 375 44. Slater, G.S.C., and Birney, E. (2005). Automated generation of heuristics for biological  
376 sequence comparison. *BMC Bioinformatics* *6*, 31.
- 377 45. Gardner, P.P., Daub, J., Tate, J.G., Nawrocki, E.P., Kolbe, D.L., Lindgreen, S.,  
378 Wilkinson, A.C., Finn, R.D., Griffiths-Jones S., Eddy, S.R., et al. (2009). Rfam:  
379 updates to the RNA families database. *Nucl. Acids Res.* *37*, D136-140.

- 380 46. Sonnhammer, E.L.L., Eddy, S.R., and Durbin, R. (1997). Pfam: a comprehensive  
381 database of protein families based on seed alignments. *Proteins* 28, 405-420.
- 382 47. Langmead, B., and Salzberg, S.L. (2012). Fast gapped-read alignment with Bowtie2.  
383 *Nature Methods* 9, 357-359.
- 384 48. DePristo, M.A., Banks, E., Poplin, R., Garimella, K.V., Maguire, J.R., Hartl, C.,  
385 Philippakis, A.A., del Angel, G., Rivas, M.A., Hanna, M., et al. (2011). A framework  
386 for variation discovery and genotyping using next-generation DNA sequencing data.  
387 *Nature Genetics* 43, 491-498.
- 388 49. Van der Auwera, G.A., Carneiro, M.O., Hartl, C., Poplin, R., Del Angel, G., Levy-  
389 Moonshine, A., Jordan, T., Shakir, K., Roazen, D., Thibault, J., et al. (2013). From  
390 FastQ data to high-confidence variant calls: the genome analysis toolkit best practices  
391 pipeline. *Curr. Protocols Bioinfo.* 43, 11.10.1-11.10.33.
- 392 50. Rozen, S., and Skaletsky, H. (2000). Primer3 on the WWW for general users and for  
393 biologist programmers. *Methods Mol. Biol.* 132, 365–386.
- 394 51. Myers, E.W., and Miller, W. (1988). Optimal alignments in linear space. *Comput.*  
395 *Appl. Biosci.* 4, 11-17.
- 396 52. Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., Marth, G.,  
397 Abecasis, G., Durbin, R., et al. (2009). The Sequence alignment/map (SAM) format  
398 and SAMtools. *Bioinformatics* 25, 2078-9.
- 399 53. Lynch, M., and Conery, J.S. (2000). The evolutionary fate and consequences of  
400 duplicated genes. *Science* 290, 1151–1155.
- 401 54. Vanneste, K., Van de Peer, Y., and Maere, S. (2013). Inference of genome duplications  
402 from age distributions revisited. *Mol. Biol. Evol.* 30, 177-190.
- 403 55. Panchy, N., Lehti-Shiu, M., and Shiu, S.-H. (2016). Evolution of gene duplication in  
404 plants. *Plant Physiol.* 171, 2294-2316.
- 405 56. Altenhoff, A.M., Škunca, N., Glover, N., Train, C.M., Sueki, A., Piližota, I., Gori, K.,  
406 Tomiczek, B., Müller, S., Redestig, H., et al. (2015). The OMA orthology database in  
407 2015: function predictions, better plant support, synteny view and other improvements.  
408 *Nucl. Acids Res.* 43, D240-D249.
- 409

## 410 **FIGURE AND TABLE LEGENDS**

411

### 412 **Figure 1. Distribution of Somatic Mutations in the Napoleon Oak**

413 **(A)** The genome of two leaf samples (outlined dots) was sequenced to identify single-  
414 nucleotide variants (SNV). 17 SNVs were confirmed and analysed in 26 other leaf samples to  
415 map their origin. A reconstructed image of the Napoleon Oak shows similar location of two  
416 SNVs (magenta dots) on the tree. Blue dots represent genotypes that are non-mutant for these  
417 SNVs. Three non-mutant samples are not visible on this projection. Location of other SNVs  
418 can be found in Fig. S1. **(B)** Location of all identified SNVs. Sectors of the tree containing  
419 each group of SNVs are represented by different colours.

420

421 **Figure 2. Napoleon Oak Apical Meristem**

422 (A) Cross-section of an apical meristem. Stem cells are delineated. Surrounding cells belong  
423 to leaf-like structures surrounding the meristem. Scale bar, 50  $\mu\text{m}$ . (B) Longitudinal section  
424 of an apical bud. Apical meristem (arrowhead) is surrounded by leaf-like structures (stars).  
425 Scale bar, 500  $\mu\text{m}$ .

426

427 **Figure 3. Spectrum of Somatic Mutations Between Two Napoleon Oak Genomes**

428 The type of substitution for 17 confirmed oak SNVs is shown.

429

430 **Table 1. *Quercus robur* Genome Statistics**

431

432 **EXPERIMENTAL PROCEDURES**

433 **Materials and Genome Sequencing**

434 Leaves were collected in April 2012 from the terminal part of a lower (sample 0) and an  
435 upper branch (sample 66) of the Napoleon Oak (*Q. robur*) on the Lausanne University  
436 Campus (Switzerland, 46°31'18.9"N 6°34'44.5"E). The age of the tree was estimated by a  
437 tree ring analysis from a sample taken at the basis of the trunk (Laboratoire Romand de  
438 Dendrochronologie, 1510 Moudon, Switzerland). DNA from the two samples was extracted  
439 and the genome sequenced. Paired-end sequencing libraries with insert size of 400 bp were  
440 constructed for each DNA sample according to the manufacturer's instructions. Then, 100 bp  
441 paired-reads were generated on Illumina HiSeq 2000 at Fasteris ([www.fasteris.com](http://www.fasteris.com)). In  
442 addition, 3 kb mate-pair libraries were constructed and sequenced. Single-molecule real-time  
443 (SMRT) sequencing (Pacific Biosciences) was performed on 25 SMRT cells according to the  
444 manufacturer's instructions (University of Lausanne Genomics Technologies Platform).

445

## 446 **Genome Assembly**

447 For sample 0, a paired-end library generated 2 x 151,194,704 reads (coverage 40X) and a  
448 mate-pair library generated 2 x 107,264,298 reads (coverage 29X). For sample 66, a paired-  
449 end library generated 2 x 158,505,474 reads (coverage 42X) and a mate-pair library  
450 generated 2 x 124,076,608 reads (coverage 33X). These reads were filtered and trimmed  
451 prior assembly using Trimmomatic (v0.3; leading:3, trailing:3, slidingwindow:4:15,  
452 minlen:36, custom adapter library) [35] and assembled using SOAPdenovo2 (v2.04.240,  
453 kmer 49) [36]. In a second step the assembly was scaffolded with mate-pairs using the same  
454 program. The assembly was further scaffolded with long single-molecule PacBio reads (25  
455 smartcells, XL-C2 and P4-C2 chemistry) and the program AHA  
456 (<http://www.pacb.com/products-and-services/analytical-software/smrt-analysis/>; SMRTPipe  
457 2.0.1 manually driven, settings (5,2,50,70), no gap-filling). Assembled sequences <1000 bp  
458 were removed to facilitate further analysis. The genome was extended with all paired-end  
459 libraries and SSPACE [37] (v2.0, -x = 1, z = 0, -k = 5, -a = 0.7, -n = 15, -T = 20, -p = 0, -o = 20, -t  
460 = 0, -m = 32, -r = 0.9) and gaps were filled using Gapfiller (v1.10, all paired-end libraries)  
461 [38].

462 We screened the paired-end libraries for potential non-oak sequences using metaphlan  
463 (v1.7.7) [39]. Based on metaphlan results, reference genomes were obtained for the non-oak  
464 genomes and the oak scaffolds were filtered against these using blast (ncbi-blast v2.28, >90%  
465 sequence identity and E-value <1e-5) [40]. The genome was next scaffolded again using the  
466 PacBio reads and PBJelly (v14.1.14) [41]. If not further specified, programs were used with  
467 their standard settings.

468

## 469 **Gene Prediction and Annotation**



470 Repetitive elements were analysed by first generating a specific repeat model using  
471 RepeatModeler (<http://www.repeatmasker.org>; v1.0.7, -engine wublast). Repetitive regions in  
472 the genome were subsequently masked with the obtained model using RepeatMasker  
473 (<http://www.repeatmasker.org>; v4.0.3). Genes were predicted by generating a *Q. robur*  
474 specific gene prediction model for Augustus (v3.0.1) [42], as described in Tran et al. [43].  
475 Instead of RNAseq reads, we used the UniProtKB reference proteome of *Glycine max*  
476 mapped with the splice aware mapper exonerate (V2.2.0, model protein2genome, geneseed  
477 250 -minintron 20, --maxintron 20000) [44]. Using this model we predicted genes and  
478 subsequently their encoded proteins for the hard-masked version of the genome (settings: no  
479 hints, no UTR predicted, no alternative transcripts). Non-coding elements were annotated  
480 using RFAM (v1.5; infernal 1.0.2; blast 2.2.26; hmmer 3.1b1) [45] in the genome with  
481 coding regions masked but repetitive elements unmasked. The predicted proteome was  
482 annotated based on homology using the FASTA toolkit  
483 (<http://www.ebi.ac.uk/Tools/sss/fasta/>; v36.3.5e) as following: proteins from the *Glycine max*  
484 proteome were first mapped with ggsearch (-b 1 -d 0 -E 1e-5 -m 8 -T 10); proteins that did  
485 not map were mapped in a next step with glsearch (-b 1 -d 0 -E 1e-5 -m 8 -T 10) and finally  
486 the rest with ssearch (-b 1 -d 0 -E 1e-5 -m 8 -T 10). The functional protein annotation was  
487 overtaken from *Glycine max*. For proteins with unknown function in *Glycine max*, we  
488 extended the annotation using the OMA database ([www.omabrowser.org](http://www.omabrowser.org)) and orthologous  
489 proteins from *A. thaliana*. PFAM [46] was used additionally to obtain functional domain  
490 annotations for the proteome and the concatenated proteome annotation was transferred onto  
491 the oak genome.

492

493 **SNV Identification**

494 First, short Illumina paired-end sequencing reads (298,735,463 and 316,299,457 for sample 0  
495 and 66, respectively) were aligned to the masked *de novo* assembly (RepeatMasker, v4.05)  
496 with Bowtie2 (v2.2.2, <https://sourceforge.net/projects/bowtie-bio/files/bowtie2/2.2.2>) [47]  
497 using default parameters. GATK [12] was used for local realignment and variant calling  
498 using standard hard filtering parameters according to GATK Best Practices recommendations  
499 [48, 49]. Prior to variant calling, each sample was screened for duplicates using PICARD  
500 tools (<http://broadinstitute.github.io/picard/>). Variants with confidence score  $\geq 50$  were  
501 retained further. We identified 1,832,554 heterozygous sites common to both samples, as  
502 well as 314,865 putative differences between sample 0 and 66 (165,489 sites predicted to be  
503 homozygous on sample 0 and heterozygous on sample 66 and 149,376 homozygous on  
504 sample 66 and heterozygous on sample 0). The distribution of the confidence scores of the  
505 1,832,554 heterozygous sites common to both samples was a superposition of a Gaussian  
506 distribution, peaking at 910, and an exponential distribution, possibly representing true  
507 positives, and the decreasing number of false positives with regard to increasing confidence  
508 score, respectively. We thus hypothesized that sites that are truly different between samples  
509 0 and 66 were unlikely to be present at sites with a confidence score below 300. We selected  
510 1,536 putative SNVs out of 314,865 for validation, on the basis of a confidence score  $\geq 300$   
511 on the heterozygous sites and  $\geq 200$  on homozygous sites. To validate putative differences  
512 between both samples, we used PCR-seq, a modification of the published RT-PCR-seq  
513 method [13]. Briefly, pairs of primers for 50-150 bp amplicons containing the targeted  
514 sequence were designed using Primer3 [50]. Touchdown PCR amplification was performed  
515 in a final volume of 12.5 ml with JumpStart REDTaq ReadyMix (Sigma-Aldrich), a primer  
516 concentration of 0.4 mM and 2 ng of gDNA per reaction in 384-well plates. Equal volumes  
517 of PCR products were pooled for each DNA template (sample 0 and 66). One ml of each pool  
518 was then purified with the QIAquick PCR Purification Kit (Qiagen) following the

519 manufacturer's instructions. The KAPA LTP Library Preparation Kit (Kapa Biosystems) was  
520 used, starting with 500 ng of purified PCR products, to create a library compatible with an  
521 Illumina sequencing platform. Clean-ups between enzymatic steps were performed with  
522 Nucleospin PCR Clean-up columns (Macherey-Nagel). After ligation of pentabase adapters,  
523 libraries were run on a 2 % agarose gel and extracted using the MinElute Gel Extraction kit  
524 (Qiagen). Libraries were sequenced on HiSeq 2000 after six cycles of amplification  
525 (Lausanne Genomic Technologies Facility). Amplicon reads were aligned, with no  
526 mismatches allowed, to a compendium of the expected amplicons that bore the reference  
527 allele, the alternate allele identified in the heterozygote sample, as well as the remaining two  
528 nucleotides at the variable position; this allowed an unbiased estimation of the error rate  
529 generated by the sequencing itself. As this method might have missed *bona fide* changes  
530 between the two sampled branches that present other heterozygous sites close by, we also  
531 aligned amplicon sequencing reads directly to the reference genome, with mismatches  
532 allowed. Only seven of the 1,536 candidates assessed were validated by PCR-seq. Sanger  
533 sequencing further confirmed them.

534         Second, to identify potential differences between both samples including masked  
535 sequences, Illumina reads of samples 0 and 66 were mapped against the non-masked oak  
536 genome assembly. The genome was 719,779,348 bp long, but 69,130,634 (9.52%) of those  
537 nucleotides were gaps and were discarded, leaving an actual search space of 650,648,714 bp.  
538 Of the latter, 458,143,725 nucleotides with a read coverage  $\geq 8$  in both samples were analysed  
539 further. The mapping process was performed at the read pair level by the genome mapping  
540 tool, fetchGWI [14], followed by a detailed sequence alignment tool, align0 [51]. Potential  
541 SNVs were called from the mapped reads by a simple read pileup process followed by  
542 detection of positions where the pileup shows variations with respect to the reference  
543 genome; this produced a list of 5,330 positions. Those positions were browsed through a

544 local adaptation of the samtools pileup browser [52] to evaluate the quality of the mapping in  
545 the surrounding region and to discriminate between well-assembled high-quality regions with  
546 two alleles per sample, or low complexity and possibly badly assembled repeated regions.  
547 Criteria for selection were  $\geq 8$  reads in each orientation (see above); 100% homozygosity site  
548 for one sample and at least 30% minor allele frequency for the other sample with variants in  
549 both orientations; and coherent sequence  $\pm 50$  bp from variant site. This manual process led to  
550 the selection of 82 putative variable positions, including the seven already identified. Upon  
551 experimental validation, 10 of the remaining 75 candidates were confirmed by PCR-seq and  
552 Sanger sequencing. The Food and Drug Administration (FDA) has evaluated this approach in  
553 an effort to assess, compare, and improve techniques used in DNA testing on human genome  
554 variation analysis (<https://precision.fda.gov/challenges/consistency>). Within this frame, our  
555 method reached a F-score (F-score evaluates precision and recall) over 95% comparable to  
556 other identifiers like BWA coupled with GATK.

557

### 558 **Estimation of the Possible Missed SNVs.**

559

560 About half of the putative variable sites with confidence scores  $\geq 200$  were assessed  
561 experimentally (1,536 out of 3,488 sites). Given the confidence scores of the tested sites, we  
562 estimated that we missed fewer than six such sites not evaluated by PCR-seq. To evaluate the  
563 number of true positives missed within candidates with confidence scores  $< 200$ , we fitted a  
564 mixture of two distributions modelled on the 1,832,554 sites that were predicted to be  
565 heterozygous on both sequenced branches (Figure S3). This suggested that we had missed  
566 fewer than eleven true SNVs. We thus estimate a total of 17 missed SNVs. Note that we did  
567 not assess the presence of larger somatic changes such as copy number variants, small indels,  
568 and transposition events.

569

570 **Estimation of the Possible False Negatives.**

571 We postulated that the false negative rate corresponds to the number of true SNVs that could  
572 have been missed in the sites that were called homozygous in both samples due to random  
573 sampling of both alleles. Assuming similar cell sizes between oak and *Arabidopsis*, there  
574 should have been ca. 133 (4,000 cm/30 cm) more mitotic divisions separating the extremes of  
575 somatic lineages in the analyzed oak, which were distant by 40 m, than in *Arabidopsis*, which  
576 reaches approximately 30 cm in height. Under the assumption that the per-generation  
577 mutation rate is correlated with the number of mitotic divisions from zygote to gametes of the  
578 next generation, the mitotic mutation rate for the oak should be between  $6.6 \times 10^{-7}$  and  $4.0 \times$   
579  $10^{-6}$  substitution/site/generation based on the estimated mutation rates in plants ( $5 \times 10^{-9}$  to  $30$   
580  $\times 10^{-9}$  substitutions/site/generation [15-17]). Since the fraction of the genome that we  
581 analyzed consisted of ca. 458 Mb with  $\geq 8$  reads in each sample, the assumption is that there  
582 should be between 305 and 1,832 SNVs. We thus calculated the probability of having missed  
583 between 305 and 1,832 true SNVs in the ca. 458 Mb (read coverage  $\geq 8$ ) of sites considered  
584 homozygous in both samples. Our selection criterion for calling a heterozygous site was that  
585 it should contain  $\geq 30\%$  of the minor allele. Given the distribution of both alleles in an  
586 heterozygous DNA sample, the probability of obtaining less than 30% of a given allele after  
587 sequencing follows a binomial distribution, decreasing as the percentage of the minor allele  
588 decreases. Thus, the added probability of obtaining less than 30% of the minor allele  
589 decreases with the number of reads. We therefore analyzed the read coverage distribution for  
590 both oak samples and sorted the data into bins. Sample 0 shows an even distribution of read  
591 coverage between 8 and 70 whereas the coverage of sample 66 ranges from 8 to 50 (Figure  
592 S4). For each coverage bin we computed the mean added probability of obtaining less than  
593 30% of the minor allele, and we multiplied this value by the fraction of the genome

594 represented in this particular bin. We next used the higher and lower values of expected  
595 SNVs (305 or 1,832) to calculate the number of SNVs possibly missed in each bin in each  
596 sample for low and high mutation rates because of random sampling of both alleles (i.e., for  
597 each bin we multiplied the added probability by either 305 or 1,832). The higher-end estimate  
598 of false-negatives for both samples was 12.9 (4.0 + 8.9) for a low mutation rate, and 79.0  
599 (24.1 + 54.9) for a high mutation rate (Table S3).

600

### 601 **SNV Genotyping**

602 Leaf DNA from different locations on the tree was prepared and amplified using primers  
603 located 100-150 bp away from the 17 confirmed SNVs (Table S4). Amplicons were then  
604 subjected to Sanger sequencing.

605

### 606 **Whole-Genome Duplication**

607 Simple clustering based on homology, (i.e., clustering the predicted proteins by identity, CD-  
608 HIT, min 90% similarity), retrieved 1,098 proteins that have a >90% identity to another  
609 protein, which is not suggestive of recent whole genome duplication. Whole genome  
610 duplication should lead to an excess of relatively old paralogs, whereas small-scale duplicates  
611 are expected to be enriched in very recent paralogs. This can be estimated from the  
612 distribution of synonymous distances (dS) [53, 54]. We computed the dS on a stringent set of  
613 4,777 paralog pairs with BLAST E-value <1e-10, removing large multigene families (more  
614 than 20 members). The distribution of dS values is clearly unimodal, with an excess of low  
615 dS values (i.e., young paralogs, Figure S5). This also does not support a recent whole genome  
616 duplication in the oak lineage.

617 To address the possibility of a more ancient duplication event, we compared our oak  
618 genome reference with itself using “BLAST all versus all” as suggested in Panchy et al. [55],

619 (i.e., similarity  $\geq 30\%$ , match length  $\geq 150$ AA and E-value  $\leq 1e-5$ ). Following this procedure  
620 we have 49,444 proteins, of which 3,650 are duplicated (7.4%), 2,070 are triplicated (4.1%)  
621 and 23.7% are present in more copies with diminishing frequency. In summary, a total of  
622 17,474 oak proteins out of 49,444 appear to be duplicated (35%), which is less than that  
623 reported for closely related species (e.g. *Medicago sativa* has about 50,000 genes of which  
624  $>75\%$  are duplicated, according to Panchy et al. [55]). We then assessed whether the  
625 similarity identified above was local, properties of similar domains, or extended along the  
626 entire protein, indicative of duplicated proteins. We found only 973 oak proteins that have  
627 duplications extending over their entire lengths. In summary, it is possible that the oak  
628 genome underwent duplication, as suggested by Panchy et al. [55], but this event appears to  
629 be rather old, as we have very few ( $<3\%$ ) duplicated genes with very high similarity ( $>90\%$ )  
630 and no second peak in the dS distribution (Figure S5). It seems unlikely that such a  
631 duplication event should compromise the identification of *bona fide* variants. Note that if the  
632 duplication would have hindered the capacity to detect these variants, they would not be  
633 found in nested sectors of the tree but rather in all 26 samples assessed.

634

### 635 **Analysis of DNA Repair Genes**

636 Orthologs between *A. thaliana*, *Prunus persica* (peach) and *Q. robur* were called using the  
637 OMA database [56]. One-to-many orthologs, e.g., between *A. thaliana* and *Q. robur*,  
638 represent duplication in the oak lineage since the divergence from *A. thaliana*; they are also  
639 known as in-paralogs of oak. We classified these in-paralogs according to whether the  
640 duplication was shared by *P. persica* and *Q. robur* (i.e., one copy in *A. thaliana* relative to  
641 several copies in both the peach and oak genomes), or whether it was peach- or oak-specific  
642 (i.e., one copy in *A. thaliana* and peach, relative to several copies in oak). The number of  
643 duplicates was reported as the number of genes that could be called duplicate (i.e., the

644 number of orthologs between each tree genome and *A. thaliana*, Table S5). We then  
645 manually compiled a list of *A. thaliana* genes involved in DNA repair from  
646 SwissProt/UniProtKB annotations (Table S6). We then counted specifically the number of  
647 duplicates for genes involved in DNA repair and reported this as the number of orthologs  
648 associated with this function (Table S2 and S7 ).  
649



650  
651  
652

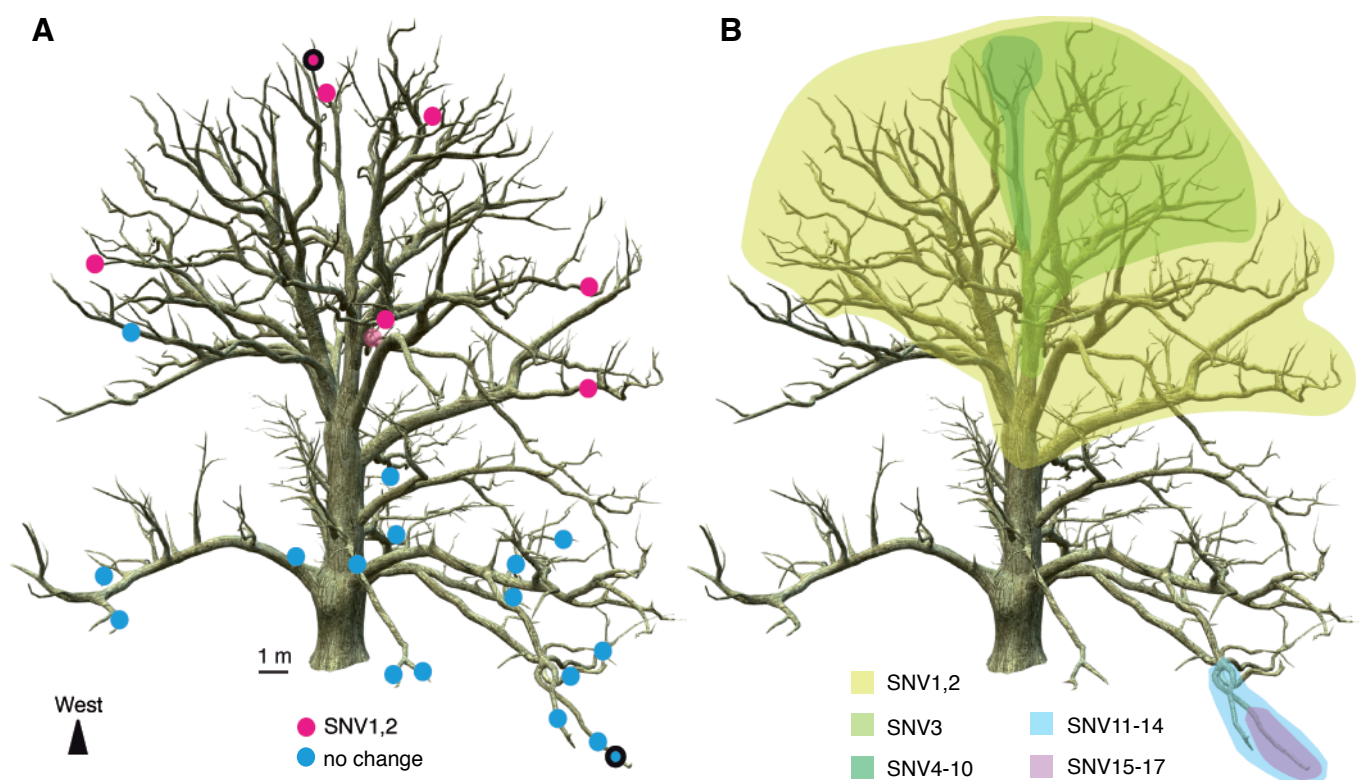
**Table 1. *Quercus robur* Genome Statistics**

---

<b>Genome</b>	
Total genome length (bp)	719,779,348
Number of scaffolds	85,557
Maximum scaffold length (bp)	317,245
NG50 based on 740 Mbp (bp)	17,014
Gaps (%)	9.52
Masked (%)	39.84
<b>Genes</b>	
Average length (bp)	2,360
Maximum length (bp)	47,221
Average intron length (bp)	740
Average exon length (bp)	232
<b>Proteome</b>	
Total predicted proteins	49,444
Full proteins	44,096
Partial proteins	5,348
Nb proteins with orthologous in <i>Glycine max</i>	39,656
Nb orthologous in <i>Glycine max</i> + functional annotation	16,323
Nb orthologous in <i>Glycine max</i> + function via ATH	23,333

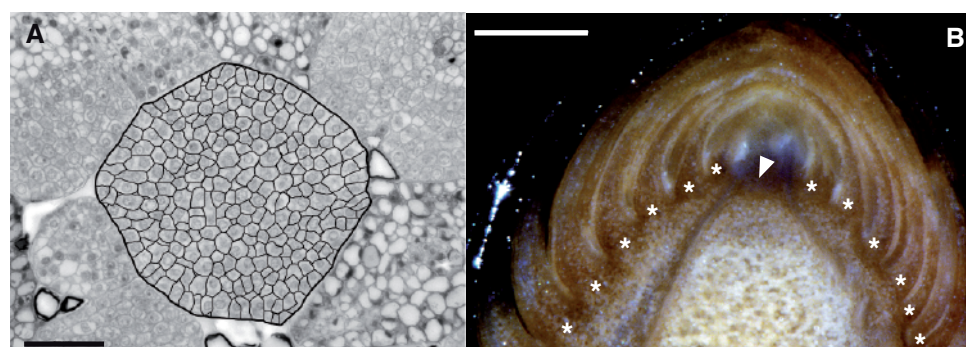
---

653



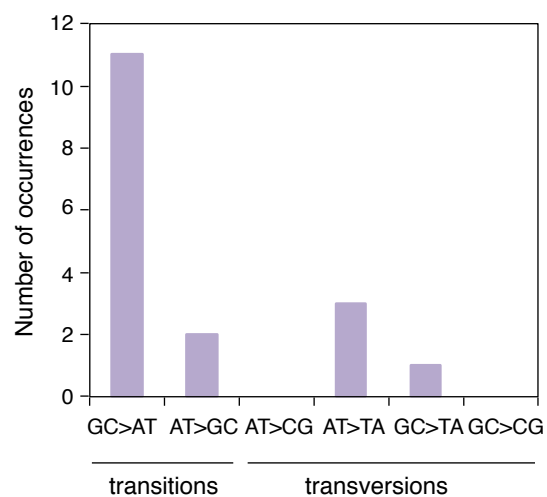
### Figure 1. Distribution of Somatic Mutations in the Napoleon Oak

(A) The genome of two leaf samples (outlined dots) was sequenced to identify single-nucleotide variants (SNV). 17 SNVs were confirmed and analysed in 26 other leaf samples to map their origin. A reconstructed image of the Napoleon Oak shows similar location of two SNVs (magenta dots) on the tree. Blue dots represent genotypes that are non-mutant for these SNVs. Three non-mutant samples are not visible on this projection. Location of other SNVs can be found in Fig. S1. (B) Location of all identified SNVs. Sectors of the tree containing each group of SNVs are represented by different colours.



## Figure 2. Napoleon Oak Apical Meristem

(A) Cross-section of an apical meristem. Stem cells are delineated. Surrounding cells belong to leaf-like structures surrounding the meristem. Scale bar, 50  $\mu\text{m}$ . (B) Longitudinal section of an apical bud. Apical meristem (arrowhead) is surrounded by leaf-like structures (stars). Scale bar, 500  $\mu\text{m}$ .



**Figure 3. Spectrum of Somatic Mutations Between Two Napoleon Oak Genomes**

The type of substitution for 17 confirmed oak SNVs is shown.