# High-quality genome and methylomes illustrate features underlying evolutionary success of oaks

Sork, Victoria L.[1,2*, **], Shawn J. Cokus[3**], Sorel T. Fitz-Gibbon[1**], Aleksey V. Zimin[4,5], Daniela Puiu[4], Jesse A. Garcia[1], Paul F. Gugger[6], Claudia L. Henriquez[1], Ying Zhen[1], Kirk E. Lohmueller[1,7], Matteo Pellegrini[3], Steven L. Salzberg[4,8]

**Affiliations:**

[1] Department of Ecology and Evolutionary Biology, University of California, Los Angeles, CA 90095-1438

[2] Institute of the Environment and Sustainability, University of California, Los Angeles, CA 90095

[3] Department of Molecular, Cell, and Developmental Biology, University of California, Los Angeles, CA 90095-7239

[4] Center for Computational Biology, Whiting School of Engineering, Johns Hopkins University, Baltimore, Maryland 21218

[5] Department of Biomedical Engineering, Johns Hopkins University, Baltimore, Maryland 21218

[6] University of Maryland Center for Environmental Science, Appalachian Laboratory, Frostburg, MD 21532

[7] Department of Human Genetics, David Geffen School of Medicine, University of California, Los Angeles, CA 90095

[8] Departments of Biomedical Engineering, Computer Science, and Biostatistics, Johns Hopkins University, Baltimore, Maryland 21218


*Corresponding author: Victoria L. Sork, vlsork@ucla.edu

**Authors contributed equally

**Short title:** Genome and methylomes of a California oak

## Abstract

The genus *Quercus*, which emerged ~55 million years ago during globally warm temperatures, diversified into ~450 species. We present a high-quality *de novo* genome assembly of a California endemic oak, *Quercus lobata*, revealing features consistent with oak evolutionary success. Effective population size remained large throughout history despite declining since the early Miocene. Analysis of 39,373 mapped protein-coding genes outlined copious duplications consistent with genetic and phenotypic diversity, both by retention of genes created during the ancient γ whole genome hexaploid duplication event and by tandem duplication within families, including the numerous resistance genes and also unexpected candidate genes for an incompatibility system involving multiple non-self-recognition genes. An additional surprising finding is that subcontext-specific patterns of DNA methylation associated with transposable elements reveal broadly-distributed heterochromatin in intergenic regions, similar to grasses (another highly successful taxon). Collectively, these features promote genetic and phenotypic variation that would facilitate adaptability to changing environments.

**Introduction**

Oaks are a speciose tree genus of the temperate forests of the northern hemisphere (from Canada to Mexico in North America, Norway to Spain in Europe, and China to Borneo in Asia) [1,2]. The genus evolved in the palearctic during a time when the earth experienced a warmer climate [3]. Fossil records indicate that sections within the genus — *Quercus*, *Lobatae*, and *Protobalanus* — were already present in the arctic during the middle Eocene 47.8–38 Mya [3]. As the planet cooled, oaks disappeared from the arctic and migrated southward, speciating as they spread over Asia, North America, and Europe. Throughout these regions, the resultant species were the foundational constituents of their plant communities [3]. This genus, which has diversified into two subgenera, eight sections, and more than 400 species [4], is an "evolutionary success story" [1]. In North America, oaks have more biomass than any other woody plant genus, including pines [5], making this genus an ecological success story as well. As dominant species, oaks play pivotal roles in shaping biodiversity, creating healthy ecosystems, and sequestering carbon needed to mitigate climate warming. Throughout human history, they have provided valuable food, housing, materials, and cultural resources across multiple continents. Here we seek insights from the oak genome to uncover mechanisms that underlie the success of oaks.

We report details of a high-quality annotated chromosome-level genome assembly for *Quercus lobata* Née (valley oak; tree SW786) and associated tissue-specific methylomes. We analyze sequence trends of heterozygosity in valley oak and the European pedunculate oak (*Q. robur*) to show that effective population size ($N_e$) has declined over time, but remained sufficiently large since divergence from a common ancestor to retain high levels of genetic variation. Large $N_e$ could help response to selection as the environment has changed over the last 50 million years. Further, our analysis of tandemly duplicated genes identifies large numbers of duplicated families, which, as Plomion et al. [6] also report, are particularly enriched for resistance genes and are likely associated with longevity and the eternal "arms race" with pests. We discover a large tandemly duplicated gene family that may be part of a previously undescribed non-self-

66   recognition system that could prevent self-fertilization and promote outcrossing, or selectively
67   allow occasional hybridizations. We also find many genes retained from the ancient γ
68   paleohexaploid duplication event of the core eudicots. These are enriched for transcription
69   factors and housekeeping genes, which may be more subject to strong (hard) selective sweeps
70   than the tandemly duplicated genes[7]. Finally, we find some surprising similarities with the
71   genomes of Poaceae (grasses — also highly successful plants). DNA methylation (BS-
72   Seq) patterns indicate heterochromatin-rich chromosome arms, and additionally show CHH
73   methylation peaks upstream of transcription start sites. Such prominent "mCHH islands" are
74   known in maize [8] and a few other plants. These features could both affect gene expression and
75   also facilitate tandem duplication events creating phenotypic variation and opportunities for
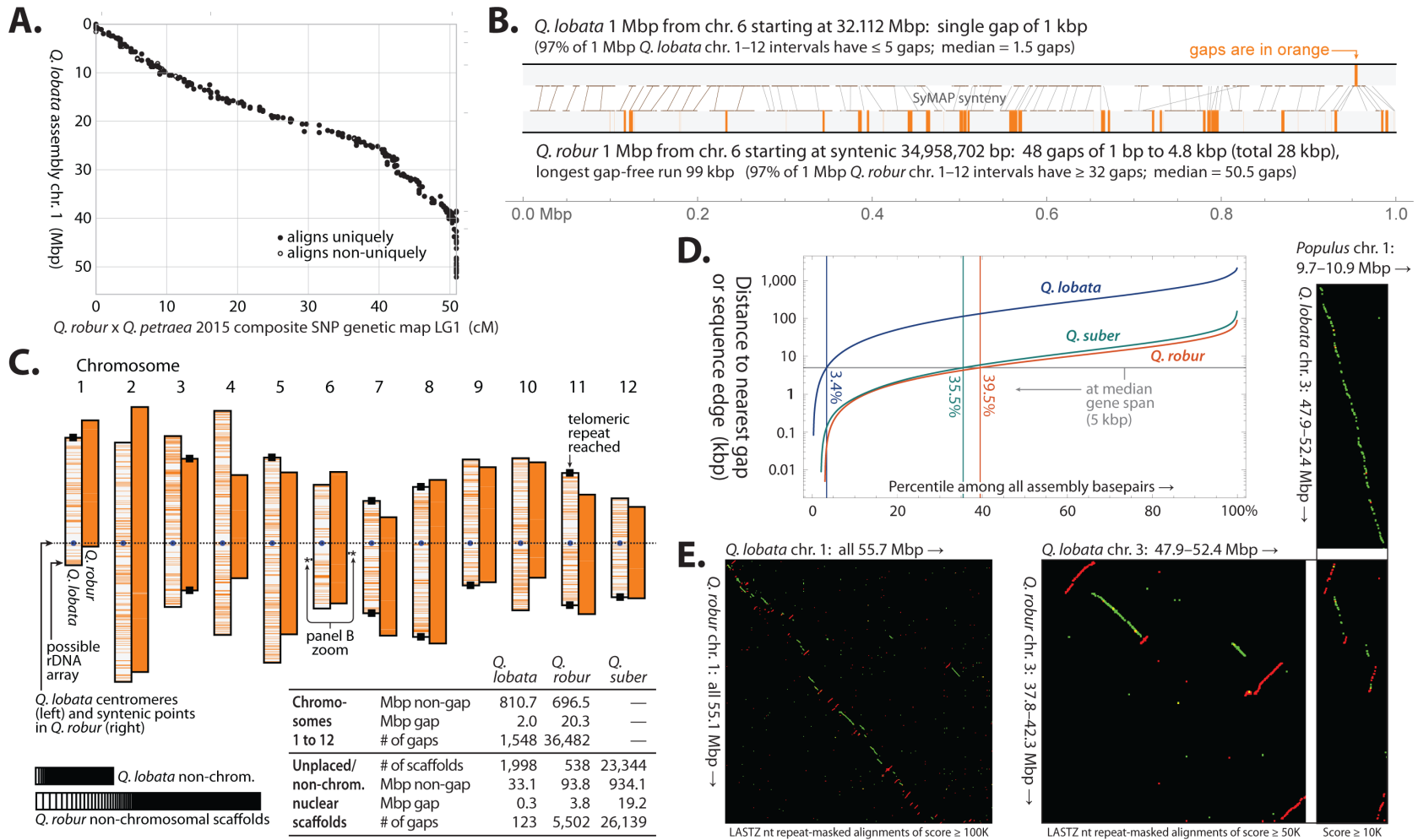76   selection.

## Results

78   *Genome assembly.* An initial draft genome (version 1.0)[9] was assembled from small (≈150x
79   coverage) and large insert (≈50x) Illumina paired end reads. The final assembly (version 3.0) was
80   constructed with the addition of Pacific Biosciences long reads (≈80x) and Hi-C long-range links
81   produced by Dovetail Genomics and the HiRise re-scaffolder[10], dramatically increasing NG50
82   scaffold size from 2 kbp to 75 Mbp (see **Methods**). The twelve longest scaffolds
83   ("chromosomes") were named and oriented to agree with pedunculate oak *Q. robur* [6], and
84   correspond (in order, but not generally orientation) with the twelve linkage groups (LGs) of an
85   existing moderate-density physical map of *Q. robur* x *Q. petraea* [11] that we did not use during
86   sequence construction. This physical map consists of 4,217 sequence context-defined single
87   nucleotide polymorphism (SNP) markers (after dropping 22 SNPs associated to two physical
88   locations each). The LGs and our chromosomes show a predominantly monotonic one-to-one
89   correspondence (e.g., chr. 1 in Figure 1A; see details in SI Section II: Validation and orientation
90   of chromosomes, Figures S1 and S2, and Table S2). 99% of SNPs had at least one BLASTN
91   alignment to our genome, and 98% of these had at least one alignment to the same
92   chromosome as its LG. (95% of SNPs had all alignments to the same chromosome, and 86% had
93   a unique alignment.) A small stretch of our chromosome 1 was found to be a mis-assembled
94   mitochondrial sequence and was replaced by a gap of the same length (Figure S3).

95   A comparison of our assembly with the two others available for *Quercus* — a chromosomal-
96   level one for *Q. robur* [6] and a short-scaffold one for cork oak *Q. suber* [12] — revealed high
97   similarity, despite ≈35M years since a common ancestor [13]. This similarity is both at the level of
98   repeats (see **Repetitive sequences**), as well as non-repetitive non-gap sequence where LASTZ
99   aligns 88% of *Q. lobata* to *Q. robur* with average nucleotide identity 96%, and 86% of *Q. lobata*
100   to *Q. suber* with average 93% identity. The larger contributing alignments tend to have even
101   higher identity; e.g., the longest alignments capturing half of *Q. lobata* have average identity
102   98% for *Q. robur* and 95% for *Q. suber*.

103   Our assembly is characterized by much higher contiguity than the other two. For example,
104   comparing *Q. lobata* vs. *Q. robur* and *Q. suber*, the number of gapless runs ("contigs") is more
105   than an order of magnitude smaller at 3.7k vs. 43k and 49k, respectively; N50 for gapless runs is

106   more than 20-fold larger at 966 kbp vs. 37 kbp and 45 kbp; and N90 is also more than 20-fold
107   larger at 205 kbp vs. 10 kbp and 9 kbp. Comparing a representative 1 Mbp from *Q. lobata* and
108   the syntenic 1 Mbp from *Q. robur* (Figure 1B), the former has a single gap of 1 kbp while the
109   latter has 28 kbp in 48 gaps of 1 bp to 5 kbp each. This pattern is typical: over all 1 Mbp regions
110   from chr. 1–12, *Q. lobata* has median 1–2 gaps (97% of regions have ≤ 5), but *Q. robur* has 50–
111   51 (97% having ≥ 32). Our assembly reaches telomeric repeats on both ends of four
112   chromosomes, and on one end of four more. (Telomeric repeats, centromeres, and rDNA are
113   discussed in SI Section V: Repetitive sequences.) Visualizing the entire *Q. lobata* and *Q. robur*
114   assemblies (Figure 1C), *Q. robur* gaps appear nearly solid. The percent of non-gap sequence
115   placed in a chromosome is 96% in valley oak vs. 88% in pedunculate oak (and 0% in cork oak).
116   The three assemblies differ considerably in total Mbp of non-gap sequence: 845 Mbp for valley
117   oak vs. 791 Mbp and 934 Mbp for pedunculate and cork oak, respectively. However, there are
118   three *Q. suber* scaffold populations by length, and the longest — those ≥ ≈50 Kbp — total a
119   more comparable 837 Mbp. More than a third of *Q. robur* and *Q. suber* base pairs are closer
120   than a median gene span (5 kbp) to an assembly gap or sequence edge, while 96% of *Q. lobata*
121   base pairs are further away (Figure 1D).

122   Apparent segmental rearrangements and inversions between *Q. lobata* and *Q. robur* were
123   unexpectedly prevalent (e.g., Figure 1E left shows chr. 1 vs. chr. 1 as typical). Most of these,
124   however, are likely scaffolding errors in *Q. robur*. Pedunculate oak has much smaller contigs,
125   and its scaffolding was constructed using linkage maps (which are low in resolution compared
126   to Hi-C) as well as synteny to *Prunus*, which may lead to mistakes in order and orientation of
127   contigs (especially for small contigs). By contrast, alignments of *Q. lobata* with more distant
128   species (*Populus*, *Eucalyptus*, *Theobroma*, and *Coffea*) showed numerous and widespread
129   regions in continuous syntenies where *Q. robur* was not as continuous; to illustrate, Figure 1E
130   right shows Mbp-scale regions of chr. 3 of the two *Quercus* vs. chr. 1 of *Populus*. Further,
131   comparison of the formerly mentioned LGs from *Q. robur* x *Q. petraea* to both the *Q. lobata*
132   and *Q. robur* assemblies shows *Q. robur* with more disagreements (Figures S1, S2). Thus, with
133   the currently available *Q. robur* assembly, we conclude that alignments of *Q. robur* versus, e.g.,
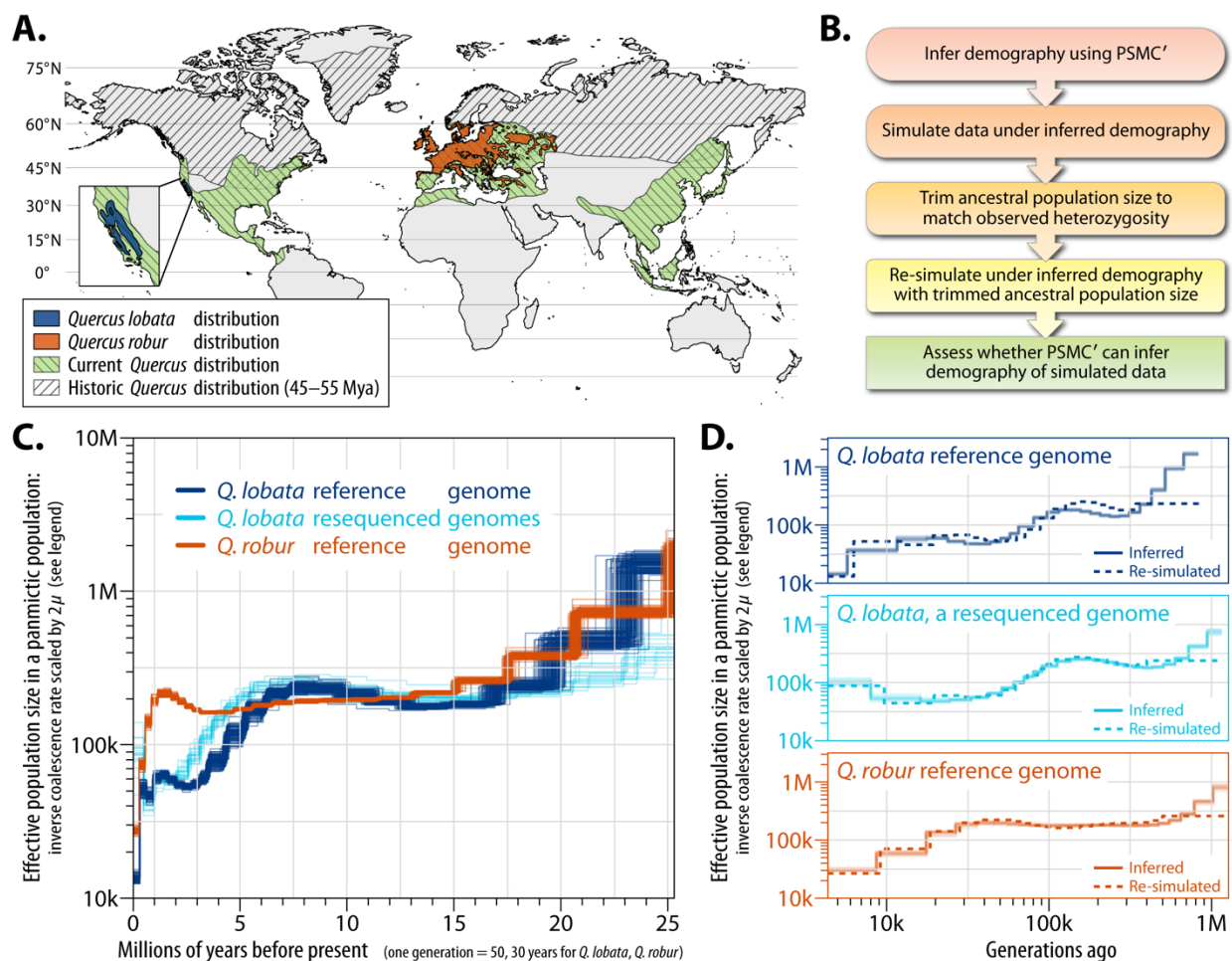134   *Q. lobata* are not reflective of true rearrangements and inversions.

**A.** *Q. lobata* assembly chr. 1 (Mbp) vs *Q. robur* x *Q. petraea* 2015 composite SNP genetic map LG1 (cM)
- aligns uniquely
- aligns non-uniquely

**B.** *Q. lobata* 1 Mbp from chr. 6 starting at 32.112 Mbp: single gap of 1 kbp
(97% of 1 Mbp *Q. lobata* chr. 1–12 intervals have ≤ 5 gaps; median = 1.5 gaps)
gaps are in orange
SyMAP synteny
*Q. robur* 1 Mbp from chr. 6 starting at syntenic 34,958,702 bp: 48 gaps of 1 bp to 4.8 kbp (total 28 kbp),
longest gap-free run 99 kbp (97% of 1 Mbp *Q. robur* chr. 1–12 intervals have ≥ 32 gaps; median = 50.5 gaps)

0.0 Mbp     0.2     0.4     0.6     0.8     1.0

**C.** Chromosome
1 2 3 4 5 6 7 8 9 10 11 12
telomeric repeat reached
panel B zoom
*Q. robur*
*Q. lobata*
possible rDNA array
*Q. lobata* centromeres (left) and syntenic points in *Q. robur* (right)
*Q. lobata* non-chrom.
*Q. robur* non-chromosomal scaffolds

|  |  | *Q. lobata* | *Q. robur* | *Q. suber* |
|---|---|---|---|---|
| **Chromosomes 1 to 12** | Mbp non-gap | 810.7 | 696.5 | — |
|  | Mbp gap | 2.0 | 20.3 | — |
|  | # of gaps | 1,548 | 36,482 | — |
| **Unplaced/ non-chrom. nuclear scaffolds** | # of scaffolds | 1,998 | 538 | 23,344 |
|  | Mbp non-gap | 33.1 | 93.8 | 934.1 |
|  | Mbp gap | 0.3 | 3.8 | 19.2 |
|  | # of gaps | 123 | 5,502 | 26,139 |

**D.** Distance to nearest gap or sequence edge (kbp) vs Percentile among all assembly basepairs →
*Q. lobata*, *Q. suber*, *Q. robur*
3.4%   35.5%   39.5%
at median gene span (5 kbp)

*Populus* chr. 1: 9.7–10.9 Mbp →
*Q. lobata* chr. 3: 47.9–52.4 Mbp →

**E.** *Q. lobata* chr. 1: all 55.7 Mbp →
*Q. robur* chr. 1: all 55.1 Mbp →
LASTZ nt repeat-masked alignments of score ≥ 100K

*Q. lobata* chr. 3: 47.9–52.4 Mbp →
*Q. robur* chr. 3: 37.8–42.3 Mbp →
LASTZ nt repeat-masked alignments of score ≥ 50K     Score ≥ 10K

**Figure 1:** Overview of assemblies of *Q. lobata* tree SW786 (version 3.0), *Q. robur* (version PM1N) [6], and *Q. suber* (version 1.0) [12]. **(A)** Alignment of a physical map linkage group 1 to *Q. lobata* chr. 1, exhibiting high concordance and overall monotonicity. **(B)** A representative 1 Mbp region from the *Q. lobata* assembly (top) and the syntenic 1 Mbp region from the *Q. robur* assembly (bottom), showing gaps in orange. **(C)** Overview of the chromosome-level assemblies (*Q. lobata* left member of each pair, *Q. robur* right) with orange lines indicating gaps, and basic statistics for all three assemblies. **(D)** Distributions of distance from a random base pair to the nearest gap or sequence edge. **(E)** Nucleotide alignments of entire chr. 1 of *Q. lobata* and *Q. robur*, showing numerous apparent rearrangements and inversions, in contrast to a more detailed illustrative region between chr. 3 of the two *Quercus* with chr. 1 of more distant *Populus trichocarpa* [14], in which the *Q. lobata* assembly is straight-line syntenic with *Populus* but that of *Q. robur* is not. Alignments between nominal same/opposite strands are colored green/red.

144    ***Demographic histories of Q. lobata and Q. robur.*** Ancient oaks evolved over 50 Mya, initially in
145    the subtropical climate of the palearctic of the Northern Hemisphere and, as the planet cooled,
146    shifting southward to their contemporary distribution throughout the Northern Hemisphere
147    (Figure 2A). Consistent with the large range, we found heterozygosity (average 0.50%–0.66%;
148    see SI Section III: Analysis of heterozygosity, and Figures S4, S5) across the genome to be similar
149    to but slightly less than the 0.73% computed for *Q. robur*, possibly due to the much larger
150    species range of pedunculate oak and/or lower representation in the *Q. robur* assembly of
151    highly homologous sequence loci resulting in increased post-alignment pileup of multiple actual
152    loci at single assembly loci. To gain insight into the population history of oaks, we inferred the
153    effective population size ($N_e$) of *Q. lobata* and *Q. robur* over time. The Pairwise Sequentially
154    Markovian Coalescent (PSMC') method [15] applied to the individuals used to build the genomes
155    mapped to their own assemblies (Figure 2B) enabled examination of the last ≈25M years of
156    evolution (Figure S7). To verify accurate inference on this timespan, we generated simulated
157    datasets using the inferred demographic history. We selected ancestral population sizes
158    matching empirical genome-wide heterozygosities (see **Methods** and SI Section IV:
159    Demographic analysis, and Figures S8, S9, and S10), and display these in Figure 2C.

160    We ran PSMC' on data simulated under trimmed demographic models and found accurate
161    inference of population size over time, except for the single oldest time step where population
162    sizes were often over-estimated (Figure 2D). The PSMC' analysis indicates ancestral populations
163    of both *Q. lobata* and *Q. robur* had high (>500k) effective population sizes that then showed
164    initially similar declines, perhaps as populations were shifting southward (Figure 2C). *Q. lobata*
165    shows an additional decline in $N_e$ at ≈5 Mya, which would have occurred after the shift from a
166    period of subtropical climate with year-round rainfall to a Mediterranean climate with summer
167    drought [16]. By contrast, for *Q. robur* (being more widely distributed throughout Europe), $N_e$
168    remained relatively flat until the last ≈1M years. At this point, *Q. robur* declines to $N_e$ < 50k (and
169    below *Q. lobata*) during the "Ice Ages" when the region was experiencing a series of warm and
170    cold periods creating genetic bottlenecks and expansions (Figure 2C and D). Both species have
171    retained sufficiently large effective population sizes to facilitate natural selection [17].

172    **Repetitive sequences.** As with many plant species, the valley oak genome contains substantial
173    repeats, with 54% of non-gap base pairs marked as repetitive by RepeatMasker in combination
174    with a species-specific database constructed by RepeatModeler+Classifier (Figure 3A and B; the
175    modeling step was essential, as RepBase only marked 13%). The largest identified portion is
176    transposable elements (TEs), primarily Copia and Gypsy elements of the long terminal repeat
177    (LTR) type. The level of repetitiveness is similar to the 54% (disregarding gaps) found by
178    application of the same process to *Q. robur* (for which Plomion et al. [6] reported 52% via REPET
179    and other annotation, including manual curation). RepeatModeler+Classifier also detects 51%
180    in *Q. suber* [12], 55% in *Eucalyptus* [18], 55% in *Theobroma* [19], 51% in *Coffea* [20] and 43% in *Populus* [14].
181    Centromeric, telomeric, and rDNA repeats for valley oak were identified (see SI Section V:
182    Repetitive sequences), and specific sequence-defined repeat superfamilies are correlated or
183    anticorrelated to various levels with centromeric proximity, forming (as do protein-coding gene
184    exons) density gradients that are the main chromosome-scale repeat-associated features,
185    presumably reflecting overall chromatin structure (Figures S11, S12, and Figure 3C–D).
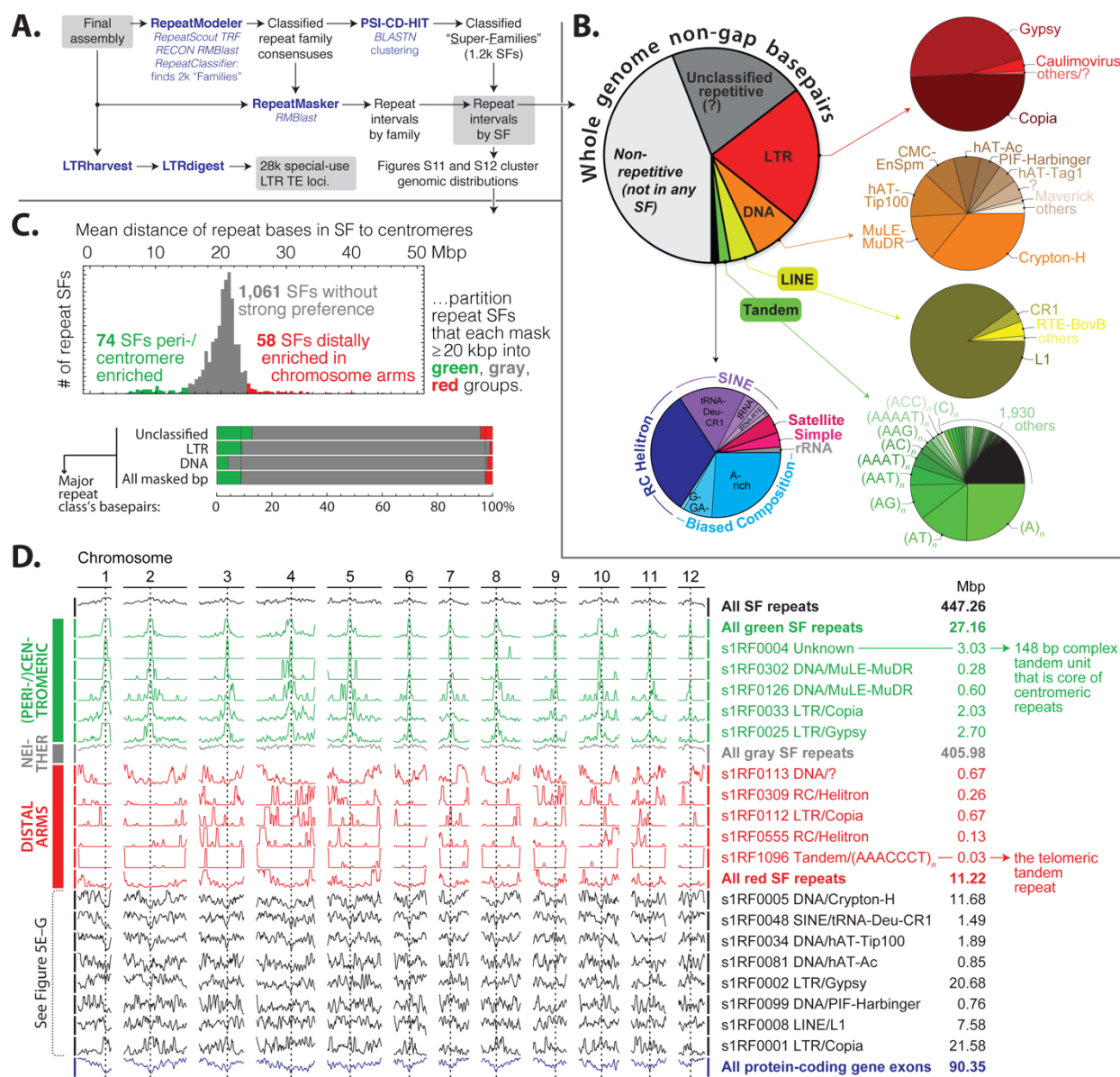
**Figure 2**: Demographic evolutionary history analysis of *Q. lobata* and *Q. robur*. **(A)** Historic and contemporary species ranges based on Barrón, et al. [3] and fossil occurrence records from the Global Biodiversity Information Facility website (GBIF.org 19th January 2019, https://doi.org/10.15468/dl.kotc15). Contemporary distribution of *Q. robur* is from the European Forest Genetic Resources Programme (http://www.euforgen.org/species/quercus-robur/); *Q. lobata* is based on Griffin and Critchfield [21]. **(B)** Stages of the analysis. **(C)** Inferred effective population sizes over time via PSMC′ (100 bootstraps shown per condition), using a mutation rate of $1.01 \times 10^{-8}$ bp per generation (see SI Section IV: Demographic analysis and Figure S6 for other parameters). **(D)** PSMC′ accurately infers demography (solid) on data simulated (dashed) under models fit to the empirical data.

The repeat content of *Q. lobata*, *Q. robur*, and *Q. suber* are very similar at the sequence level. There are six combinations in which RepeatModeler can be used to build a species-specific repeat consensus database from one of the three *Quercus* assemblies, which can then be applied by RepeatMasker to one of the two other assemblies. In all six combinations, 89% to 92% of non-gap base pairs are marked the same way (repetitive or not repetitive) as when the native consensus database for the species being masked is used.

**Figure 3:** Dispersed and local (tandem/satellite, simple/biased composition) repetitive sequence in *Q. lobata*. **(A)** Primary analysis outline. **(B)** Assembly partitioned into RepeatClassifier/RepeatMasker major and minor classes; 54% of non-gap base pairs are covered by repeat superfamilies (SFs), and transposable elements (TEs) are prevalent. **(C)** Unsupervised comparisons of how the 1,193 individual SFs each with ≥ 20 kbp distribute across chromosomes (Figures S11 and S12) suggest the primary distributional diversity at chromosome scale is proximity to centromeres (green, 74 SFs totaling 27 Mbp) vs. telomeres (red, 58 SFs totaling 11 Mbp) vs. more-or-less uniformity (gray, 1,061 SFs totaling 406 Mbp). **(D)** Chromosomal distribution of selected SFs and sets of SFs, illustrating the diversity across and within the trichotomy of (C). The *y*-axis in each row is linear number of member base pairs in 3 Mbp bins every 1 Mbp, with zero at the lower edge and 95th percentile (or row maximum if the percentile is zero) at the upper edge. Black rows near the bottom are the representative SFs of Figure 5E–G.

**Gene prediction and annotation**. Using the AUGUSTUS gene modeler [22] and a diverse set of experimental data (Iso-Seq, RNA-Seq, DNA methylation) and *in silico* data (known proteins, repeats), we modeled 68k putative protein-coding genes (PCGs) (see **Methods**, Figure 7A and

215   Table S3). As many corresponded to transposons with little expression or appeared
216   hypothetical for other reasons, we removed 29k to obtain the primary set of 39,373 PCGs we
217   report, of which 35k have at least one intron and all of which have UTRs annotated and are
218   ostensibly complete. *Q. robur* reports only 29k PCG models, of which just 20k have introns, and
219   about half UTRs; in the other direction, *Q. suber*'s annotation by NCBI (thinned to one isoform
220   per locus) reports more 49k PCG loci (about half with UTRs), but a more comparable 36k with
221   introns and 38k ostensibly complete, and with a much higher number containing transposon
222   domains by comparison.

223   We assigned gene names, functions, and orthologs via the PANTHER and Pfam components of
224   InterProScan, and OMA [23]. We evaluated the *Q. lobata*, *Q. robur*, and *Q. suber* scaffolds and
225   single isoform PCG model sets with BUSCO (Figure 7B). *Q. lobata* compares favorably to the
226   other two, and does not have the high multicopy anomaly of *Q. suber* in the 303-USCO ODB9
227   Eukaryota set [24], or the high missing and fragmented fraction of *Q. robur*'s small model set
228   (especially with the more comprehensive 2,121-USCO set for Eudicotyledons from ODB10).

229   **Gene duplications.** Protein–protein alignments among the *Q. lobata* PCGs exposed a rich
230   panoply of duplication structure in terms of genomic positions, ages, and functions. Prominent
231   and complex tandem-like blocks of high-similarity genes can be seen via visualizations of all–
232   vs.–all alignments (see **Methods**). These duplications often involve local rearrangements, and
233   can extend into megabases with dozens of genes involved at a time. Figure 4A (left third)
234   exhibits two illustrative 5 Mbp regions of chr. 4. Approximately 40% of PCGs participate in these
235   blocks, which have sizes of two to ≈100 genes each, with larger sizes rarified like a power law
236   (Figure S13). Roughly a third of participating genes are duplicated only once, slightly more than
237   half two to 20 times, and only a tenth more than 20 times. Visualizations (e.g., coordinated
238   Figure 4A middle third) of the synonymous codon substitution rate ($K_s$) over gene pairs in blocks
239   suggest a wide variety of ages for the majority of retained expansion for individual blocks.
240   Larger blocks tend to be older (Figure 4F colored distributions), but even old blocks tend to
241   have younger points suggestive of ongoing growth. While numerous tandem gene copies are
242   shorter or have reduced or no RNA-Seq evidence of expression, many copies (even within larger
243   blocks) are not particularly short or of lower expression and so do not appear to be
244   pseudogenes. Functions of tandemly duplicated genes are diverse, as evident from the variety
245   of Pfam domains they contain (e.g., coordinated Figure 4A right third). Relatively few distinct
246   domains, however, are strongly enriched over all tandemly duplicated genes, and include NB-
247   ARC, LRR_8, B_lectin, LRR_1, TIR_2, LRRNT_2, p450, TIR, and PGG (associated with
248   resistance/defense); Pkinase_Tyr and Pkinase (signal transduction); UDPGT (the large UDP-
249   glucoronosyl/glucosyl transferase family); S_locus_glycop, PAN_2, and DUF247 (see below); F-
250   box, FBA_3, and FBA_1 (protein–protein interactions/degradation, signal transduction and
251   regulation); and GST_N, GST_N_3, and GST_N_2 (glutathione S-transferases, with functions
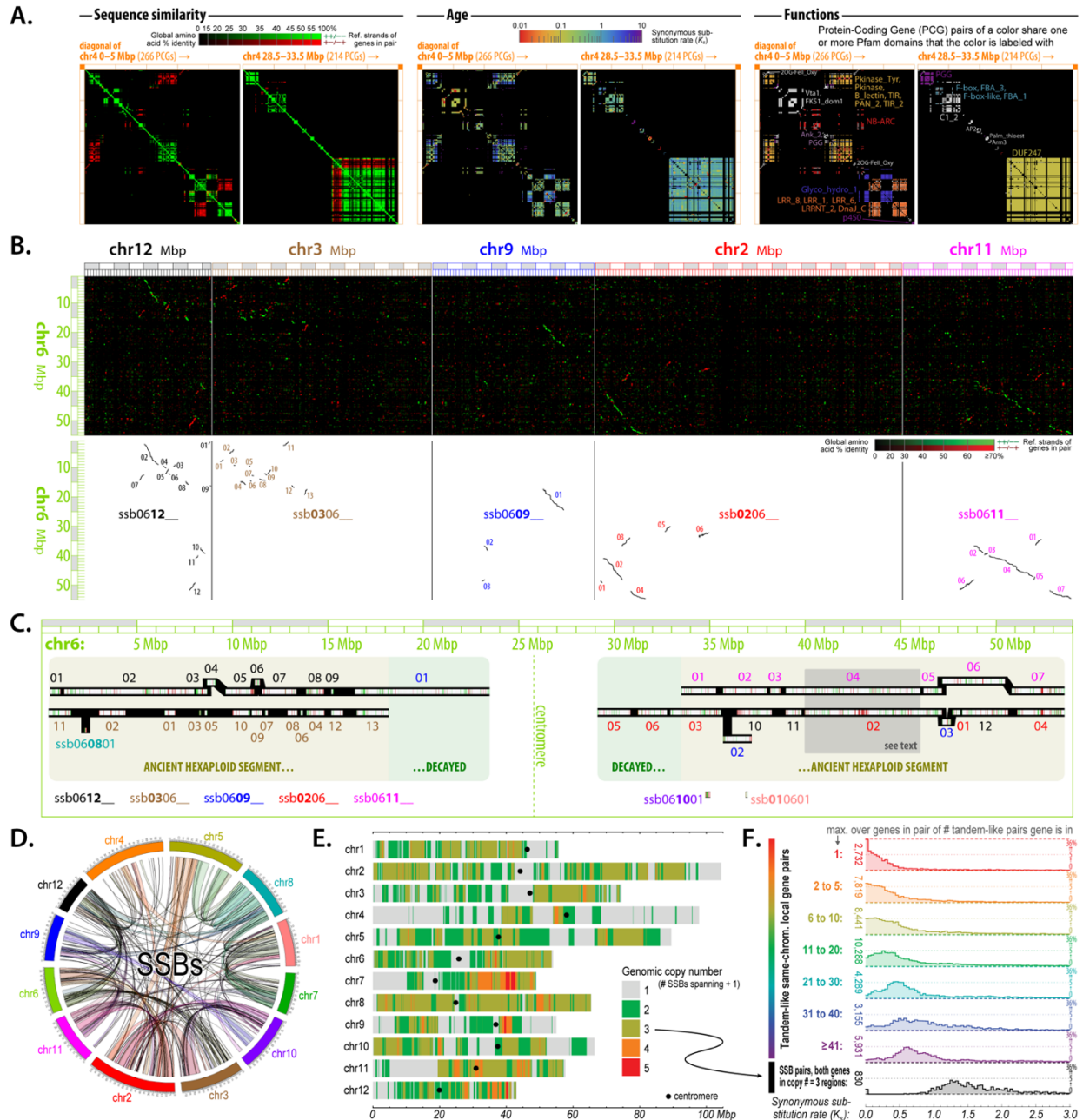252   including stress tolerance/signaling and detoxification).

253   Many of the strongly enriched domains are part of the domain architecture of plant disease
254   resistance genes (R-genes) identified by Gururani, et al. [25].  It is difficult to be sure whether an
255   R-5gene is actually acting as a pathogen defense mechanism in a given plant species, but

256    Gururani, et al.[25] reviewed the experimenal evidence and identified eight classes of R-genes
257    based on the arrangements of domains and structural motifs. Using their criteria for domain
258    combinations ( see Supplementary Information, Section VI. R-gene identification), we analyzed
259    the domain architecture of 39373 *Q. lobata* proteins and found 751 R-genes, which contained
260    the highly likely combination of domains involved in R-genes, and another 2176 genes that are
261    good candidates for R-genes because they include enough qualifying features (Table S4).  For
262    the *Q. robur* annotation (25,808 proteins) [6], we counted 632 R-genes plus 1645 candidate R-
263    genes and for *Q. suber* (49,388 proteins) [12], we found 723 *R*-genes plus 2182 candidate R-genes
264    (Table S4). These numbers are based on the predicted gene models from each genome rather
265    than DNA sequences, and the differences among species are more likely to represent
266    differences across annotations rather than differences in DNA sequences.  Collectively they
267    document high levels of R-genes in oaks and illustrate tremendous opportunity for plant
268    defense mechanisms.

269    ***Possible DUF247-based non-self-recognition system.*** An investigation of PCGs found in blocks
270    containing at least 30 tandemly duplicated genes uncovered DUF247 (PF03140) as the most
271    enriched Pfam domain (Table S5; also see large block in Figure 4A). The only known suggested
272    function for DUF247-containing genes ("DUF247 genes") is from the Poaceae family, where two
273    DUF247 genes in rye grass segregate with each of two known self-recognition loci and are
274    proposed to be the male determinants of a multi-locus self-incompatibility system[26,27]. Among
275    the evidence is a self-compatible rye grass species with a disrupted DUF247 gene [26]. If the
276    DUF247 gene family affects self-recognition in oaks, the extensive duplication suggests a non-
277    self rather than a self-recognition system [28,29]. This type of system has been demonstrated in
278    Solanaceae, e.g., petunia [30] and tomato [31], Rosaceae, e.g., pear and apple [32], and Plantaginaceae,
279    e.g., snapdragon [33]. In these, the S-locus includes a single female determinant gene (S-RNAse)
280    and commonly seven to 20 linked paralogs of male determinant F-box genes (SLFs). In
281    snapdragon [33], up to 37 linked male determinant SLF genes were observed, while (at the other
282    extreme) *Prunus* species have a single F-box gene for the male determinant and appear to have
283    adopted an S-RNAse-based self-recognition system rather than non-self-recognition [34],
284    demonstrating the feasibility of transitioning from one to the other. The span of the large
285    DUF247 block (Figure 4A) contains 34 predicted PCGs with a complete DUF247 domain, 22 with
286    partial DUF247 domains, and 17 additional genes. Among the 17 additional genes are two
287    pectinesterase inhibitor-like genes shown to be involved in regulating pollen tube growth in
288    maize [35], a pectin depolymerase gene, two E3 ubiquitin ligases that have been shown to confer
289    self-incompatibility when transplanted to *Arabidopsis* [36], a DNA helicase, and 11
290    uncharacterized genes. DUF247 genes are entirely specific to plants and usually carry a single
291    copy of the domain that comprises almost the entire gene. Across the 104 plant genomes in
292    Pfam Release 33.1[37] with $\geq$ 17,500 predicted protein entries (to restrict to genomes most likely
293    to be complete), the top five by number of DUF247 domain occurrences are three tree species
294    — *Juglans regia* (English walnut) $n$ = 201; *Eucalyptus grandis*, $n$ = 188; and *Populus trichocarpa*
295    (black cottonwood), $n$ = 161 — and two polyploid cultivars (wheat, $n$ = 192, and peanut, $n$ =
296    165). These tree species, like *Q. lobata* ($n$ = 186), do not have identified incompatibility systems,
297    are frequently highly outcrossing, and sometimes self-fertilize at low levels.

298    **Long-surviving duplicated genes.** Also striking in the visualizations of protein alignments were
299    self-syntenic blocks (SSBs): syntenic runs of proteins within *Q. lobata*, generally between
300    different chromosomes, with a variety of lengths and gene pair densities. Figure 4B (top) shows
301    chr. 6 vs. chr. 12/3/9/2/11 as exemplary (although in low resolution per limited space). For
302    further analysis, 236 SSB runs, each with four to hundreds of gene pairs, were extracted (e.g.,
303    Figure 4B bottom) and given accessions "ssbXXYYZZ" with XX ≤ YY indicating the chromosomes
304    involved and ZZ as serial number; more than 7,100 PCGs are directly involved. High resolution
305    examination made evident that, on any given chromosome, runs tended to end and begin close
306    by, and for any particular point on a chromosome to be covered by very few runs (typically,
307    zero to two), so that (nearly) disjoint SSBs could often be clearly ordered to form a small
308    number of chromosome-scale chains (Figure 4C black bars). While a few recent segmental
309    duplications appear, most SSBs are likely "ghosts" of the ancient genome triplication polyploidy
310    event γ associated with early diversification of the core eudicots, thought to have occurred
311    about 120 Mya [38-40]. The high age of many SSBs is supported by the synonymous substitution
312    rate ($K_s$) for gene pairs in SSBs in triplicated regions being very high (almost entirely > 1.0;
313    Figure 4F black distribution), as well as the generally short length and scattered nature of SSBs
314    (which are within *Q. lobata*) compared to syntenies between *Q. lobata* and different species
315    (*Populus*, *Eucalyptus*, *Theobroma*, and *Coffea*).

316    While general triplication is clear from the gene pair-defined SSBs (e.g., Figure 4C white bars,
317    with green and red showing supporting gene pairs), few syntenic gene triples have been
318    retained, and detection and characterization of the γ ghosts would be unrepresentative for an
319    analysis restricted to gene triples. For example, in the gray shaded region of Figure 4C involving
320    chr. 6/2/11 and spanning 320 chr. 6 genes, the 59 chr. 6 genes supporting local one-to-one
321    chr. 6/2 synteny have only eleven chr. 6 genes in common with the 39 supporting local one-to-
322    one chr. 6/11 synteny. Even before chaining as in Figure 4C, two thirds of the genome are in a
323    SSB (Figure 4D and E), with the largest fraction (34%) actually covered by two (consistent with
324    triplication) and 27% by one (decayed triplications and a few recent segmental duplications); a
325    third (34%) is not covered, and only 5% is covered by three or four SSBs (likely duplications post
326    triplication). Relative to all genes, those in one or two gene pairs supporting SSBs tend to be of
327    higher expression with lower repetitive sequence in their immediate vicinity, and are enriched
328    for certain functional classes, including transcription factors and housekeeping genes
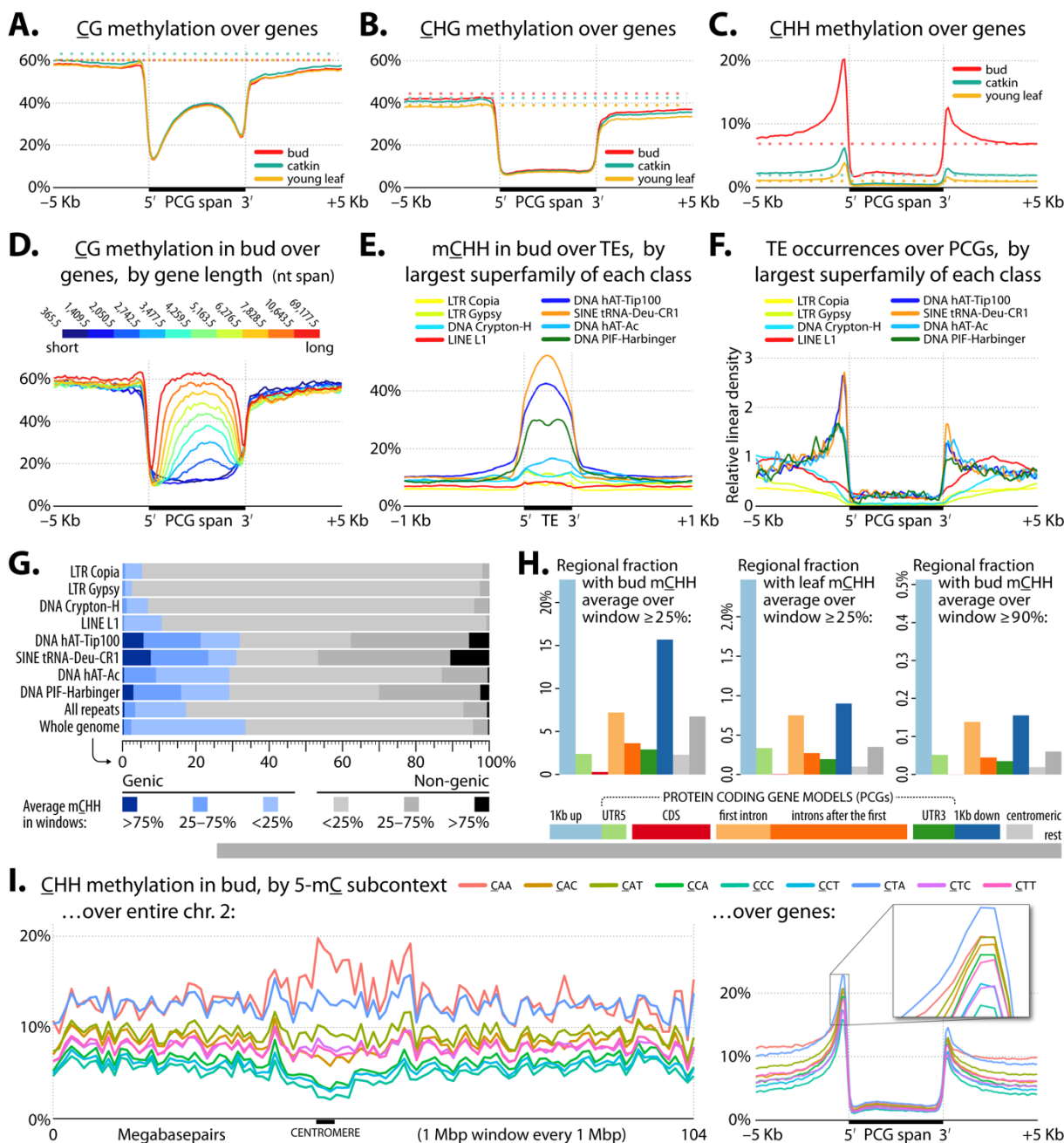329    (Table S6).

**Figure 4:** Duplicated protein-coding genes. **(A)** Sequence similarity (amino acid identity), age (synonymous substitution rate $K_s$), and functions (shared Pfam domains) for all pairs of proteins within two illustrative 5 Mbp regions of chr. 4. Nearly half of *Q. lobata* PCGs are involved in tandem-like blocks of varying sizes (up to Mbp scales and dozens of genes at a time), often locally rearranged, and originating and growing at a variety of ages. Genes involved are diverse, but enriched in certain functions. **(B, C)** With no recent whole-genome polyploidization, most of the detected PCG syntenies of *Q. lobata* to itself ("SSBs") are small and diffuse and reflect the core eudicot triplication event γ over 100 Mya. Despite its age, this event remains quite evident — albeit highly fragmented, dispersed, and partially decayed. The whole of chr. 6 vs. the whole of chr. 12/3/9/2/11 are shown as exemplary. **(D, E)** SSBs [even without chaining as in (C)] cover much of the chromosomes. The highest fraction (34% of base pairs) is spanned by manifest triplication, 27% by duplication (while some duplication is recent, most appears to be decayed triplication), and 34% by no detected extant synteny. **(F)** The pairwise synonymous substitution rate ($K_s$) tends to be very low for-genes tandemly duplicated just once (red) and increases as tandem-like block size increases (orange to violet), suggesting larger blocks are older. $K_s$ is essentially always extremely high (≥ ~1.0) for SSB gene pairs where both pair-genes lie in chromosomal regions spanned by exactly two SSBs (black), supporting the syntenic triplications to be of ancient origin.

345   ***Genome-wide patterns of DNA methylation and strong mCHH islands***. Whole-genome bisulfite
346   sequencing for bud, catkin, and leaf tissue revealed mean DNA 5-methylcytosine methylation
347   (BS-Seq) levels in CG (mCG) and CHG (mCHG) nucleotide contexts as relatively stable across
348   tissues (Figure 5A and B), while levels in CHH (mCHH; Figure 5C) were notably higher in bud
349   than catkin and young leaf, likely due to the increased proportion of undifferentiated meristem
350   tissue [41]. Mean levels for regions surrounding genes are similar to genome-wide means for all
351   tissues in all contexts (mCHH 1–7%, mCHG 39–43%, mCG 60–62%; Figure S14), with the
352   exception of peaks of mCHH near transcription boundaries of genes (Figure 5C). These mCHH
353   peaks are similar in both position and scale above background to the "mCHH islands" of maize
354   [42,43] (Auxiliary Spreadsheet 1). We examined mCHH across representative repeat superfamilies
355   (SFs), specifically, those of highest mass within selected RepeatClassifier minor repeat classes,
356   as seen in Figure 5E. Within genic regions, three SFs — s1RF0048 ("SINE tRNA-Deu-CR1"),
357   s1RF0034 ("DNA transposon hAT-Tip"), and s1RF0099 ("DNA transposon PIF-Harbinger") —
358   were both high in mCHH and preferentially located in the highly methylated gene boundary
359   regions (Figures 5E and F). Members of these SFs are found in both genic and non-genic regions
360   with broadly similar mCHH levels (Figure 5G and Figure S15). However, in view of overall
361   genome-wide mCHH levels (including centromeres and intergenic space), we find regions
362   surrounding genes to be highly enriched for mCHH (Figure 5H). Similar enrichment patterns are
363   seen in bud and leaf, despite different overall mCHH levels (Figure 5C and H), and similar
364   patterns are also seen if mCHH window stringency is varied from 25% to 90%, although at these
365   extremes we observe decreases in the relative amount of downstream and non-genic mCHH
366   (Figure 5H). All methylation is typically low near transcription boundaries (Figure 5A), and
367   remains low for mCHG and mCHH across gene bodies. However, gene body mCG rises for
368   longer-genes, reaching near-background levels in the longest genes (Figure 5D).

369   ***Broad distribution of heterochromatin.*** *Q. lobata* appears to have heterochromatin dispersed
370   throughout chromosomes more or less equally, with only minor increase of density toward
371   centromeres. This interpretation is based on both the distribution of genes and repeats as well
372   as indications of widespread histone-driven DNA methylation, a pattern more similar to maize
373   and rice methylomes than to the *Arabidopsis* and tomato methylomes in which the methylated
374   repeats are concentrated in pericentromeric heterochromatic regions [44,45]. As such, a majority
375   of repeat mass does not show strong positional correlation with centromeres (Figure 3D gray).
376   Also, 92% of PCGs have a RepeatMasker-defined repeat within the gene's upstream 2 kbp,
377   which is high, because among 34 angiosperms, reported numbers range from 29% (*Arabidopsis*)
378   to 94% (*Zea mays*), with an average of 50% [43]. (See also Auxiliary Spreadsheet 1).
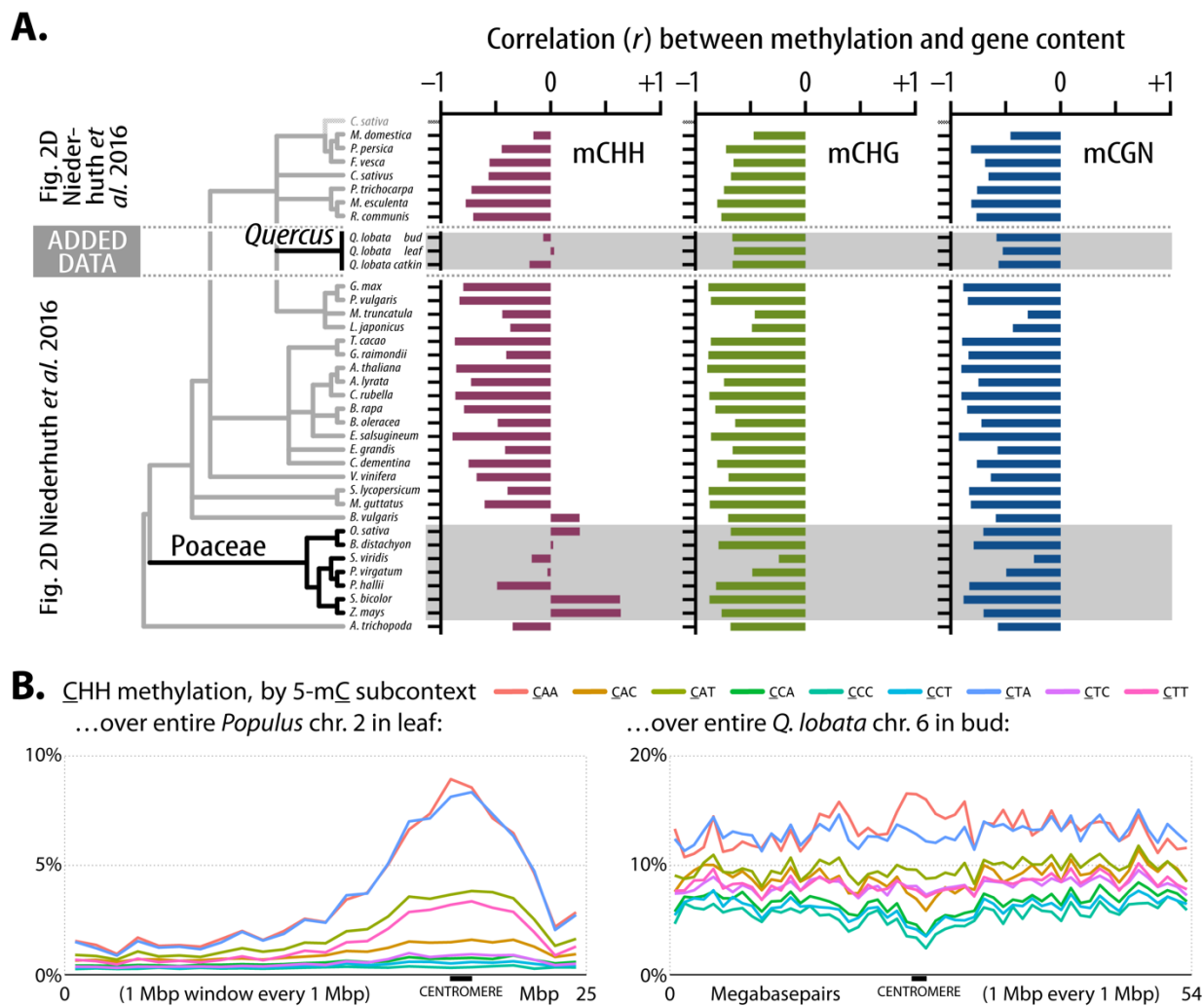
**Figure 5:** *Q. lobata* DNA methylation in protein-coding genes and repeats. **(A–C)** Average methylation levels (100 bp windows) with respect to PCGs (normalized to 5 kbp long) for the three sampled tissues (bud, catkin, and young leaf) by methylation context: **(A)** CG, **(B)** CHG, and **(C)** CHH. Dotted lines show genome-wide backgrounds, and TSS/TES = Transcription Start/End Site. **(D)** mCG for-genes in deciles by gene length. **(E)** Average bud mCHH (20 bp windows) across representative repeat SFs (normalized to 400 bp long) in selected RepeatClassifier minor classes. **(F)** Relative density of representative repeat SFs around genes (100 bp windows). **(G)** Distribution of mCHH for representative repeat SFs (100 bp sliding disjoint windows). 'Genic' = gene spans enlarged by 1 kbp on each end. **(H)** Partitioning of whole genome base pairs into nine types of regions vs. mCHH coverage. Lower horizontal bars reflect relative size. Vertical bars show percent of each genomic context covered by 100 bp windows with mCHH > 25% or > 90% in bud or young leaf. **(I)** mCHH by 3 nt subcontext (queried cytosine is underlined, and is the first nt of the three); l**eft side:** 1 Mbp windows across all of chr. 2, **right side:** across genes (normalized to 5 kbp length) in bud tissue.

392   A second indication of heterochromatin-rich chromosome arms is the type of methylation
393   found on intergenic repeats.  Different mechanisms of generating plant mCHH, such as RNA-
394   directed DNA Methylation (RdDM) or CMT3-mediated histone-associated methylation, have
395   been shown to have distinct preferences for specific nucleotide subcontexts (finer than
396   CG/CHG/CHH: CAA vs. CAC vs. …). Histone-associated mechanisms are typically responsible for
397   methylation of heterochromatin and have much stronger biases than RdDM [44]. *Q. lobata* has
398   strong CHH subcontext preferences at chromosome scale (Figure 5I left and Figure S16). Bias
399   patterns around centromeres are likely to indicate the general methylation pattern of
400   heterochromatin in oaks, while chromosome arms represent a mix of genes and intergenic
401   spaces. The peaks of mCHH surrounding gene boundaries (i.e., the mCHH islands) show a
402   distinct pattern, with preference for CAA strongly reduced (Figure 5I right). Moving from gene
403   boundaries toward intergenic space, the subcontext pattern progressively reverts to the likely
404   heterochromatic signal of the centromeres (Figure S17).

405   An additional measure of the similarity in genome organization between oaks and grasses is the
406   level of correlation between methylation and gene count across chromosomes. When we
407   augment Figure 2D from Niederhuth et al. [43] with oak findings, oak is again found comparable to
408   the grasses (Poaceae) and less to the other studied angiosperms (Figure 6A). The low mCHH
409   and gene count correlation reflects a combination of unusually strong gene boundary mCHH
410   islands relative to the background mCHH level (Figure S18 and Auxiliary Spreadsheet 1) and low
411   average gene density in chromosome arms (Figure S19).

**Figure 6:** (**A**) Pearson correlation (*r*) between methylation level and number of genes (100 kbp windows) for mCHH (left), mCHG (middle), and mCG (right) context levels from leaf tissue-based Figure 2D of Niederhuth et al. [43], inserting our values for three oak tissues (bud, leaf, and catkin from tree SW786, having matched analysis details as closely as possible). (**B**) Comparison of all nine DNA subcontext methylation levels within the CHH context over an illustrative chromosome of *Popular trichocarpa* [46] and *Q. lobata*. (See Figure 5 legend.)

418    Discussion

419    Our analysis of a high-contiguity, chromosome-level annotated oak genome reveals previously
420    unreported features of oaks that might contribute to its ability for adaptation to new
421    environments and resulting dominance in North American ecosystems. We find surprising
422    similarities to grasses (Poaceae), another highly successful group of plants. Oaks and grasses
423    both have genomes with large repeat-rich intergenic regions and share methylation features
424    that are somewhat unusual, given the current sampling of methylomes in the literature.
425    Interest has been growing in the adaptive potential provided by large complex intergenic
426    regions often  found in plants with larger genomes [47-49]. For example, a substantially higher
427    percentage of loci associated with phenotypic variation are found in the large intergenic regions
428    of maize versus the smaller intergenic regions of *Arabidopsis* [49]. Much of this regulatory
429    variation has been found in non-TE stably unmethylated DNA [50,51], such that more than 40% of
430    phenotypic variation in maize was associated with open chromatin that makes up less than 1%
431    of the genome [52]. On the other hand, high density of diverse TEs, which has been connected
432    with local adaptation [53], can be a source of both transcription factor binding sites and regulatory
433    non-coding RNAs [54], and play a role in three-dimensional genome structure [51,55,56]. An
434    abundance of intergenic heterochromatin-like structure has been demonstrated in grasses[8,57,58]
435    and, based on patterns suggestive of histone-driven methylation [44,45], are likely also found in
436    oaks (Figure 5I, Figures S16, S17, S19 and Auxiliary Spreadsheet 1). Given the dramatic
437    differences reflected in the chromosome-wide subcontext methylation patterns in the gene-
438    dense arms of *Arabidopsis* and tomato versus the wider spread of genes in maize and rice[44],
439    and similar differences in poplar versus oak (Figure 6B and Figure S20), oaks and grasses may
440    have some regulatory strategies distinct from those in other angiosperms. Another indicator of
441    similarity between oaks and grasses is the correlation of CHH methylation levels (mCHH) and
442    gene count along chromosomes (Figure 6A). A comprehensive characterization within oaks and
443    across the angiosperms awaits further experimentation and better, more comparable genome
444    sequences, constructed and annotated with consistent methods.

445    Pronounced mCHH islands are another feature shared between oaks and grasses. In maize,
446    mCHH islands have been proposed to enforce boundaries between heterochromatin and
447    euchromatin, and as such contribute to maintaining suppression of TEs during increases in
448    neighboring gene expression [8,42,59]. Measured as the ratio of peak mCHH to whole genome
449    average mCHH, we find oaks have unusually strong 5' mCHH islands (Auxiliary Spreadsheet 1),
450    but it remains to be seen if they also contribute to boundary enforcement. It is possible they
451    are simply the result of the type of TEs found near gene boundaries. In valley oak (Figure 5),
452    maize [60], and *Arabidopsis* [61], mCHH is influenced by TE family, proximity to genes, and
453    chromosomal location. The strong enrichment of small, highly methylated TE families near-
454    genes (Figure 5E and F) could be due to, for example, selection against large TEs in gene
455    proximal regions.

456    A potentially exciting discovery is the presence of many Pfam DUF247 domains in one of the
457    largest and densest blocks of tandemly duplicated genes (Figure 4A), as these domains could be
458    part of a non-self-recognition compatibility system [29]. DUF247 genes have been implicated in a

459  self-recognition system of ryegrass [26,27], analogous to S-RNAse-based self- and non-self-
460  recognition systems in petunia [30,31], and tomato, apple, snapdragon, and peach [62]. Oaks have
461  long been thought to possess some kind of self-incompatibility system because of their high
462  outcrossing rates, but the single gene SI systems have not fit observations. However, a non-self-
463  recognition system would be consistent with observed crossing results among self, intra-, and
464  inter-specific pollinations [63]. Both self- and non-self-recognition systems of co-adapted genes
465  expressed in pollen and pistil and preventing self-fertilization have evolved independently in
466  several lineages of angiosperms [29,64]. While the roles of DUF247-containing genes need
467  experimental verification, their large numbers and high diversity at the amino acid level are
468  consistent with a non-self-recognition system that could both promote outcrossing while also
469  permitting occasional self and interspecific crosses.

470  Oaks have a vast reservoir of tandemly duplicated genes of a wide variety of ages (Figure 4F),
471  contributing to their genetic and phenotypic diversity. As reported for pedunculate oak [6],
472  resistance genes are a particularly prominent component of the tandemly duplicated gene
473  blocks in valley oak, especially the larger (and older) ones (Auxiliary Spreadsheet 2: see
474  worksheets for tandem pairs >20, 30, 40).  The three oak genomes contain hundreds to
475  thousands of potential R-genes: 732 to 2927 for *Q. lobata*, 632 – 2247 for *Q. robur*, and 793 to
476  2905 for *Q. suber*.  In defending oaks from bacteria, viruses, nematodes, oomycetes, and
477  insects, these R-genes may both enable the long lifespan of oaks[6], and also address the puzzle
478  of how a single or two oak species are able to dominate so many of the ecosystems they
479  occupy. The classic Janzen–Connell ecological hypothesis proposes that pathogens promote
480  tropical forest diversity through conspecific negative density-depending (CNDD) mortality, but
481  CNDD has been shown across all forest types [65,66]. In oaks, the high number and potential
482  complexity of R-genes could provide a mechanism to reduce CNDD mortality caused by
483  pathogens [67]. Moreover, the large effective population size could maintain R-genes, especially if
484  not costly [68]. In fact, other ecosystem-dominant trees, which also contain large numbers of
485  domains associated with resistance genes (such as, NB-ARC and LRRs), include the highly
486  speciose *Eucalyptus* (~600 species) and *Populus* (Table S7). Extensive research demonstrating
487  the importance of R-gene diversity at both the individual and the population level is ongoing in
488  *Arabidopsis*, crop species and other plants [69,70].  Studying oaks swith large and complex pools of
489  R-genes will provide an important extension of this work.

490  Inspection of our highly contiguous genome identified numerous syntenic blocks of remnant
491  genes from the γ triplication event, which occurred ≈120 Mya ago when the common ancestor
492  of angiosperms underwent two whole genome duplication events [38-40]. More than 18% of
493  protein-coding genes participate in a gene pair directly supporting a self-syntenic block (SSB),
494  and more than a third of the genome is spanned by a manifest triplication (even without
495  chaining blocks). SSBs (for example, Figure 4 and Auxiliary Spreadsheet 2) provide an extensive
496  single genome resource for documenting remnants associated with the γ event. Our annotation
497  finds triplicate families to be enriched for transcription factors, as well as signal transduction
498  and housekeeping genes generally (Table S6 and Auxiliary Spreadsheet 2), as has been found in
499  other studies, e.g., Rensing [71]. These genes, although maintained over millions of years and
500  highly interconnected [72], can respond to selective pressures modifying their existing roles. For

501  example, a recent study of silver birch found selective sweeps around candidate genes enriched
502  among ancient polyploid duplicates that encode developmental timing and physiological cross-
503  talk functions[7]. In oaks, it would be constructive to learn whether these ancient genes have
504  undergone positive selection, allowing adaptation to new environments.

505  Genomes of high-quality document the deep evolutionary history of species. The oak genome
506  has many features that provide hints of possible reasons for their success. Our exploration has
507  uncovered several surprising similarities to the highly diverse grass genomes that may indicate
508  analogous or even homologous adaptive strategies that would increase functional diversity in
509  addition to the diversity generated by extensive gene duplications. Future oak studies may
510  benefit by looking to the extensive experimental results from both wild and crop grasses for
511  clues to potential mechanisms contributing to their evolutionary success.

## Methods

513  ***Study species, samples, and genomic lab work****. Quercus lobata* Née (Fagaceae) is a widely-distributed endemic
514  California oak species found in oak savannas, oak woodlands, and riparian forests. Oak have a highly outcrossed
515  mating system[73] with the potential for long distance gene flow occurring through wind-dispersed pollen with long-
516  tailed distributions, despite many near-neighbor pollinations [74,75]. Acorn dispersal is often restricted except for
517  occasional long-distance colonization by jays [76]. Occupying an unglaciated region of California, contemporary
518  populations are at least 200k years old with no evidence of severe bottlenecks during cold periods [77,78] like those
519  described for the European oaks from glaciation that retreated in the last 10k–20k years, allowing rapid
520  recolonization from refugia in Italy and Spain [79]. Valley oak and other California oak species have been used as a
521  reliable food source and cultural resource by native peoples of western North America for the last 10k years [80].
522  Since the arrival of Europeans, valley oak populations have experienced extensive habitat loss [81], and current
523  population recruitment is jeopardized by cattle grazing, rodents, and other factors [82,83]. Moreover, as its climate
524  niche shifts north and upward [82,84,85], extant populations are further challenged by climate warming.

525  The focal tree for this study is *Q. lobata* adult SW786, which is located at the UC Santa Barbara Sedgwick Nature
526  Reserve, is the same individual that was sequenced for version 1.0 of the genome [9]. Leaf samples for the initial
527  Illumina sequencing (532M paired-end [PE] 250 nt reads with ≈500 nt inserts giving 133 Gnt and ≈175x coverage,
528  and 318M mate pair [MP] 150 nt reads from ≈3 knt to ≈12 knt fragments giving 48 Gnt and ≈56x coverage) were
529  collected in September 2014, as described in Sork, et al. [9]. Additional leaves were collected and DNA extracted in
530  April 2016 for Pacific Biosciences whole genome SMRTbell libraries (6M reads of mean ≈9 knt and N50 ≈13 knt
531  giving 58 Gnt and ≈80x coverage), and in March 2017 for Dovetail Chicago Hi-C library preparation. For details of
532  the 19 whole genome resequencing libraries (Illumina PE, mean ≈24x coverage) used for the demographic analysis,
533  three-tissue (bud, leaf, stem) Pacific Biosciences Iso-Seq and Illumina RNA-Seq transcriptome libraries contributing
534  to annotation, and three-tissue (bud, catkin, and young leaf) whole-genome bisulfite libraries (Illumina SE, ≈18x –
535  19x coverage) for the DNA methylomes, (see SI Section I: Sample collection, library preparation, sequencing, and
536  initial data processing).

537  ***Genome assembly.*** We constructed the final genome in multiple stages. Stage 1: For the initial "Hybrid Primary"
538  assembly (818 Mbp in 3.6k scaffolds, with longest 6.7 Mbp and NG50 ≈1.2 Mbp assuming at-the-time estimated
539  730 Mbp for the haploid genome), we applied MaSuRCA 3.2.1 [86] to our genomic Illumina PE, Illumina MP, and
540  PacBio SMRT reads. The assembler identified high heterozygosity and selected diploid settings, allowing it to set
541  aside most divergent haplotype variants; the result generally contains a single haplotype, but randomly phased, as
542  we chose the larger scaffold whenever the assembler split two haplotypes into distinct scaffolds. Those scaffolds
543  filtered out as alternative haplotypes were gathered into the "Hybrid Alternative" additions (466 Mbp in 17k
544  scaffolds, with longest 1.2 Mbp). Stage 2: To assist completeness, we aligned to Stage 1 Primary+Alternative 82k of
545  84k transcripts and gene fragments from a prior RNA-Seq-derived transcriptome [87], with 81k aligning to Primary.

546    To avoid loss of potential coding regions, we moved 317 scaffolds from Alternative to Primary, forming the
547    "Hybrid-plus-Transcript Primary" assembly (872 Mbp in 4.0k scaffolds, with longest 6.7 Mbp and NG50 ≈1.2 Mbp),
548    and "Hybrid-plus-Transcript Alternative" additions (412 Mbp in 16k scaffolds, with longest 0.8 Mbp). Stage 3: We
549    increased NG50 by aligning Stage 2 Primary scaffold ends with bwa mem [88], merging scaffolds that had unique end
550    matches of > 94% identity longer than 40 kbp. This created the "Hybrid-plus-Transcript-Merged Primary" assembly
551    (861 Mbp in 3.2k scaffolds, with longest 10.2 Mbp and NG50 ≈1.9 Mbp) and "Hybrid-plus-Transcript-Merged
552    Alternative" additions (16k scaffolds). Stage 4: Next, we generated Hi-C long-range linking information from an
553    Illumina-sequenced library produced by Dovetail Genomics, which we used to re-scaffold with HiRise [10] after read
554    alignment with a modified SNAP (http://snap.cs.berkeley.edu), dramatically increasing NG50. Scores from the
555    HiRise learned likelihood model were used to identify and break presumed misjoins, identify prospective joins, and
556    commit joins above a threshold; shotgun reads from Stage 1 were used to close gaps where possible. Stage 5:
557    Finally, after HiRise, any redundant haplotype contigs remaining (that truly belong in the same place as the other
558    haplotype in a scaffolded assembly) are expected to be adjacent to the other haplotype as this is as close as they
559    can be placed under the linear ordering constraint of HiRise output. We used this property to remove remaining
560    extra haplotype contigs by aligning adjacent contigs to each other and finding those smaller than their direct
561    neighbor that had > 50% syntenic alignment with the neighbor, thereby moving 14 Mbp to Alternative and forming
562    the final "Hi-C-Scaffolded-plus-Neighbor-Cleaned Primary" ("version 3.0") assembly (Figure 1C). The twelve longest
563    scaffolds represent near full-length chromosomes (Figures S1 and S2) and total 811 Mbp (96%) of non-gap
564    sequence.

565    ***Comparisons of Q. lobata and Q. robur assemblies to linkage map.*** The *Q. robur* x *Q. petraea* linkage groups
566    (LGs)[11] are taken from http://arachne.pierroton.inra.fr/cgi-bin/cmap/map_set_info?map_set_acc=51 using Table
567    S3 from Lepoittevin, et al. [89] as sequence-defined SNPs, dropping SNPs associated to more than one LG. Genomic
568    locations were identified with BLASTN+ 2.2.30 ($E < 10^{-15}$, word size 8), keeping for each query all alignments with
569    bitscore ≥ 97% of the top bitscore. We plot SNPs that have either a unique surviving alignment, or multiple
570    alignments but all to the same chromosome and with chromosomal span of hits ≤ 2 Mbp wide.
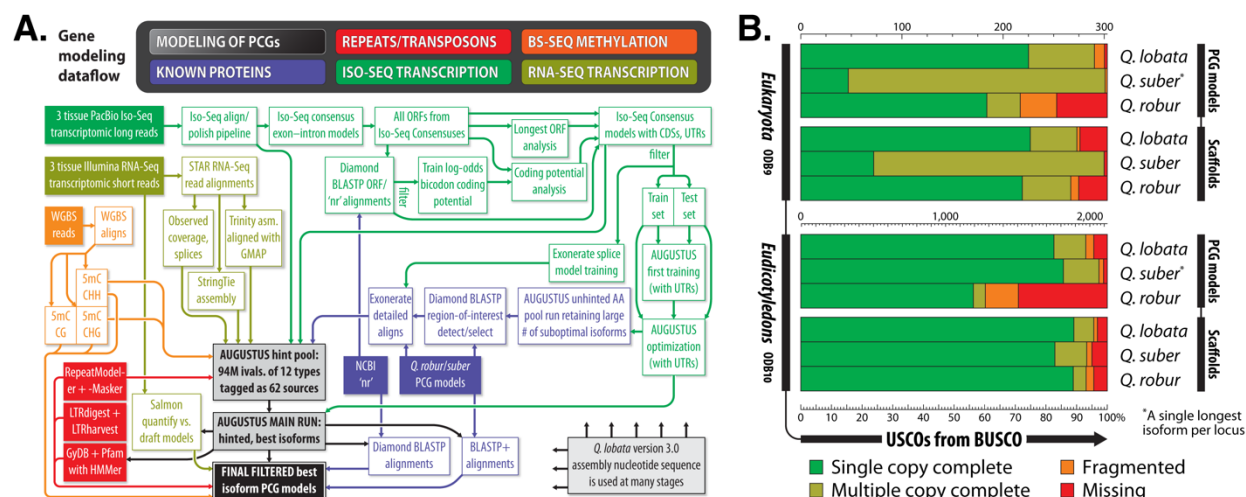
571    ***Nucleotide- and amino acid-level alignments and 2-D visualizations of sequence similarity, $K_s$, and shared***
572    ***Pfam domains, variously within and between genomes for Q. lobata, Q. robur, Q. suber, Populus, Eucalyptus,***
573    ***Theobroma, and Coffea.*** Various alignments and visualizations appear in Figures 1E, 4A, and 4B, and at the project
574    website (TBD) and assisted in *Q. lobata* with discovery and identification of the tandem-like blocks of duplicated
575    genes and syntenic self-syntenies (SSB). The principal software components were LASTZ for nucleotide alignments
576    (with masked repeats from RepeatMasker after RepeatModeler); BLASTP, Diamond, and Parasail for homologous
577    gene pair detection (on respective genome project protein-coding gene models) and subsequent detailed
578    alignment refinement; C++, Mathematica, and Perl for scripting and pixel generation/import; and ImageMagick
579    and Adobe Photoshop for manipulation, browsing, and annotation of generally multi-gigapixel images of high
580    resolution (e.g., 10 kbp/pixel). Pfam hits were determined with InterProScan or direct HMMer runs. $K_s$ was
581    computed for all homologous protein pairs discovered with other tools by re-aligning with 'needle' from EMBOSS
582    (http://emboss.sourceforget.net/), converting to the level of codons with 'pal2nal.pl'
583    (http://www.bork.embl.de/pal2nal/), and finally computing $K_s$ with 'codeml' from PAML
584    (http://abacus.gene.ucl.ac.uk/software/paml.html).

585    ***Demographic history.*** We inferred demographic history using the PSMC' algorithm [15] by mapping the *Q. lobata* and
586    *Q. robur* sequencing project genomic shotgun reads to their respective reference assemblies, as well as high-
587    coverage genomic reads for 19 *Q. lobata* individuals to the *Q. lobata* assembly (Table S1). We called heterozygous
588    sites in each genome (forming a VCF file with all callable sites) and composed input for PSMC' with
589    vcfAllSiteParser.py (https://github.com/stschiff/msmc-tools). Masking and filters are as described in SI Section IV:
590    Demographic analysis — Input to PSMC'. We ran PSMC' using default parameters except 200 for maximum number
591    of iterations. Because PSMC' inference can be prone to biases, we assessed robustness of our conclusions
592    (Figure 2B). To determine if inference is affected by re-use of the same reads as used to build the reference
593    assembly, we analyzed the 19 re-sequenced *Q. lobata* individuals beyond the reference individual SW786. These
594    showed similar population size changes (Figure 2C light blue) as SW786 (Figure 2C dark blue), suggesting little bias
595    from re-use. We also assessed if PSMC' is capable of accurate inference by generating simulated datasets following

596 the inferred demographic history, and re-ran inference on these (Figure 2D). These runs suggested that the only
597 major issue was the oldest population sizes often being over-estimated. We thus selected ancestral population
598 sizes matching empirical genome-wide heterozygosities (Figures S8, S9, and S10) and trimmed display in Figure 2C
599 accordingly (Figure S7 exhibits untrimmed trajectories). Finally, we tested whether PSMC' could reliably infer
600 changes in population size on timescales relevant to *Quercus*. We simulated 10 test datasets of each run type
601 under our presented demographic models (in Figure 2C) using the coalescent simulator msprime [90]. With each
602 simulated genome, we computed heterozygosity and used PSMC' to infer demography; see SI Section IV:
603 Demographic analysis — Simulations in msprime. We found accurate inference of population sizes over time,
604 except for the single oldest time step where it tends to be over-estimated (Figure 2D). Note that inferred
605 demographic trajectories from whole genome-based methods such as PSMC' can be complex but not predict
606 empirical summary statistics such as the genome-wide distribution of heterozygosity [91].

607 ***Repetitive sequences***. The primary repeat analysis is outlined in Figure 3A, and began with construction of a
608 *Q. lobata*-specific database of repeat families by RepeatModeler/Classifier open-1.0.8, which was then applied
609 with RepeatMasker open-4.0.6. Family consensus sequences are not always full length for their class or
610 irredundant by close sequence similarity; we applied PSI-CD-HIT 4.7 to family consensus sequences at 45%
611 nucleotide identity (the level where, as the threshold was lowered, intracluster similarities stopped falling in
612 frequency and began rising) and chose a canonical rotation and strand for tandem repeat units, so as to cluster
613 families into repeat "superfamilies" (SFs). Generally, each SF was assigned the RepeatClassifier class of the longest
614 member of the family that was not unknown (if any; approximately two-thirds of SF-covered base pairs were
615 classifiable). Annotated intervals for a SF are the nucleotide-level union of all intervals for member families, and
616 SFs were assigned "s1RF####" accessions roughly serialized by descending mass. For certain uses (e.g., gene
617 annotation), we also applied structurally-aware LTRharvest and LTRdigest from GenomeTools 1.5.9 to specifically
618 target the abundant LTR TEs, identifying 28k instances of total mass 184 Mbp (not much larger than the 179 Mbp
619 in LTR-classified SFs). Further details are in SI Section V: Repetitive sequences.

620 ***Annotation of protein-coding genes.*** Figure 7A outlines dataflow of the PCG modeling process employed. ***Pure Iso-***
621 ***Seq models.*** The Pacific Biosciences of California, Inc. (PacBio) pipeline generated 197k–223k nominally full length
622 non-chimeric polished transcripts ("reads") from the poly-A-selected strand-specific bud, leaf, and stem PacBio-
623 sequenced Iso-Seq libraries. Pooling tissues, Minimap2 aligned each read to zero to five reference genome
624 locations (96%–99% uniquely); 11% of alignments were filtered out based on empirical criteria. Preliminary exon–
625 intron structure was obtained by focusing on the reference side of each alignment, ignoring short insertions and
626 merging short deletions and gapless blocks. Inspection found compact and well-isolated gene loci with generally
627 concordant pileups at each of the 24k tentative loci, which had highly variable coverage (1 to 14k reads each; 56%
628 ≥5, 22% =1). Inspection of higher-coverage loci found reads within pileups to vary: (i) at the exact coordinate level
629 (with exon–intron boundaries moved by typically < 10 nt vs. common); (ii) at the structural level (introns resized or
630 deleted or inserted, generally in a small minority of reads, and at more loci and in more ways than likely by
631 alternative splicing); and (iii) in extent (with some reads truncated, especially at the 5' end, with loss of multiple
632 exons possible). A consensus exon–intron model per locus was generated by resolving (i) via rounding boundaries
633 within ±25 nt to a most common boundary; (ii) by generally keeping exons and introns only in at least half of reads;
634 and (iii) by extending 5' and 3' ends to the furthest extent observed. CDS assignments were made considering
635 three methods: (i) longest ORFs; (ii) filtered (including restriction to only near-best hits per locus) BLASTP-
636 equivalent (Diamond 0.9.22.123) alignments ($E < 0.001$) of translations of all ORFs to the entire NCBI 2018-05-18
637 'nr' database (22k consensuses had at least one hit, with 83% of top hits involving at least half of both the
638 translated ORF and the NCBI sequence, and with 99% having ≥ 50% amino acid identity, and 90% having $E < 10^{-35}$;
639 assignment of an ORF required agreement among all surviving hits); and (iii) a $\log_2$-odds bicodon coding potential
640 trained using a selected subset of the NCBI analysis. Partial (and six-frame) ORFs were permitted. A consensus was
641 assigned CDS (i.e., an ORF) if (ii) identified an ORF, the longest member of (i) was of the same frame with non-
642 empty intersection with that ORF, and (iii) was also of the same frame and with non-empty intersection. This
643 attributed CDS (and, hence, UTR5 and UTR3) to 19k loci, with 95% on the consensus read strand and 85% having ≥
644 50 nt of both UTR5 and UTR3. Hand inspection of a random subset found them to be of generally good quality
645 (often needing no edits). Diverse AUGUSTUS hints were constructed from the Iso-Seq reads and pure Iso-Seq
646 models for eventual use in the final AUGUSTUS run near the end of the PCG modeling process.

647

**Figure 7: (A)** Dataflow of protein-coding gene (PCG) modeling. **(B)** BUSCO v3 analysis of PCG models and genomic scaffolds for *Q. lobata*, *Q. robur*, and *Q. suber* against the ODB9 Eukaryota and ODB10 Eudicotyledons USCO sets.

**AUGUSTUS bootstrap.** The 19k pure Iso-Seq models were filtered to a very high confidence subset of 2,639, then thinned to 2,558 by choosing single representatives from homology clusters determined via Exonerate 2.4.0 affine:local protein alignments. These were split uniformly at random into 1,698-model training and 860-model test sets and used to bootstrap AUGUSTUS via "new_species.pl" (enabling UTRs) and "etraining", then optimized with "optimize_augustus.pl", and also used to train the splice model of Exonerate. **RNA-Seq.** Paired end 101+101 nt Illumina HiSeq 4000 RNA-Seq reads were also collected from rRNA-depleted strand-specific bud, leaf, and stem libraries. The 121M to 153M pairs per tissue were aligned with STAR 2.5.3a and assembled into nominal transcripts with reference guidance by StringTie 1.3.4d and Trinity 2.6.6; Trinity output was aligned back to the reference genome with GMAP 2017-11-15. A large collection of diverse AUGUSTUS hints were constructed from STAR observed genomic base pair coverage and empirical splices, StringTie reference-quoted transcripts, and GMAP alignments. **Known proteins.** Protein translations of the *Q. robur* and *Q. suber* PCG models were BLASTP-equivalent (Diamond) aligned to a temporary trained/optimized but unhinted AUGUSTUS run generating and retaining a very large number of suboptimal isoform models. The resultant hits were used to identify regions of interest on the *Q. lobata* reference genome, that were then aligned vs. *Q. robur* and *Q. suber* in splice-discovery detail with splice model-trained Exonerate. Numerous strong AUGUSTUS hints were then constructed from Exonerate's alignments. **Repeats.** Reference genome base pairs masked by RepeatMasker from the species-specific database constructed by RepeatModeler were weakly hinted to AUGUSTUS as non-exonic. **DNA mCHG and mCHH (BS-Seq) patterns.** Similar to repeats, sufficiently high mCHG or mCHH levels from merging the three tissues of the DNA methylation analyses were weakly hinted to AUGUSTUS as non-exonic. (Both *a priori* expectation and empirical examination of preliminary AUGUSTUS runs without methylation-based hinting had these marks as very highly anti-correlated with PCGs. mCG was not used, as it is complex, being high both in repeats and in many PCGs due to gene body methylation of non-short genes.)

**Main AUGUSTUS run:** The above data sources provided 94M hinting intervals of 12 types tagged as from 62 sources. (As AUGUSTUS scores cannot be configured to be continuous functions of hint evidence strength [e.g., numeric coverage level from RNA-Seq], continuous strengths were generally broken into small numbers of discrete bins, with fixed scoring per bin.) A three-line patch (in extrinsicinfo.cc) to the AUGUSTUS C++ source code was required to enlarge hard-coded limits. One top isoform model per locus was predicted by the trained, optimized, UTR-aware, and now hinted main AUGUSTUS run. **Filtered final PCG models.** Numerous models from the main AUGUSTUS run were, e.g., clearly transposons with no or little evidence of observed expression. Based on several indicators (including Salmon-quantified per-model RNA-Seq expression, overlap with annotated repeats, presence of LTRdigest/harvest or GyDB/HMMer transposon domains, average mCG and mCHG and mCHH levels, and

681   Diamond and BLASTP+ alignments with NCBI 'nr' and *Q. suber* and *Q. robur* PCGs), we removed such and other
682   hypothetical models with poor evidence.

683   ***Enrichment analyses.*** Benjamini-Hochberg false discovery rate (FDR)-adjusted hypergeometric *p*-values were used
684   to determine Pfam do main enrichment in targeted subsets of tandemly duplicated genes and genes within SSBs.

685   ***Methylomes and analysis of tissue-specific methylation patterns.*** Sample collection, library preparation,
686   sequencing, and initial methylation calling are described in SI Section VII: Methylomes and analysis of methylation
687   patterns. Libraries were prepared using the TruSeq Nano DNA (Illumina) and Epitek kits (Qiagen), and sequenced
688   as 100 nt single end reads on an Illumina HiSeq 4000 to median coverage 18–19-fold. Methylation levels were
689   determined using Methylpy v1.4.6 [92]. DeepTools v3.1.2 [93] computeMatrix and plotProfile were used to assess
690   methylation levels with respect to gene models and repeat superfamilies (Figures 5A–F, Figure 5I, Figures S15 and
691   S17, with default parameters except as described in the legend). Methylation levels for 100 bp windows were
692   calculated by dividing the total number of reads calling 'T' (= methylated) by the total number of informative reads
693   ('C' or 'T') for all genomic cytosine positions in the appropriate sequence context within the window. Genome wide
694   average methylation levels (Figure 5 and Figure S14) were calculated by averaging 100 bp window levels for the
695   twelve chromosomal scaffolds. Per-site methylation levels in Figure S14 were calculated by dividing reads showing
696   methylation ('T') by all informative reads ('C' or 'T') for each position, plotting with R ggplot2 v3.3.2 [94]. Designation
697   of genomic regions with respect to genes (1 kbp up, 5' UTR, etc.) was done with Bedtools v2.27.1 (BEDTools: a
698   flexible suite of utilities for comparing genomic features https://doi.org/10.1093/bioinformatics/btq033) and
699   bed12toAnnotation.awk (https://github.com/guigolab/geneid/blob/master/scripts/bed12toAnnotation.awk). PCG
700   model spans do not overlap in our annotation; however, overlaps for 1 kbp upstream and 1 kbp downstream
701   regions were removed from the 1k up and 1k down categories, including overlaps that spanned neighboring genes.
702   Gene regions overlapping with intervals (200 kbp to 3 Mbp) covering peri-centromere regions were removed.
703   Introns were separated into first intron vs. other introns. Chromosome scale plots of subcontext methylation
704   (Figure 5I and Figure 6A) were calculated with Bedtools as the mean of the percent methylation at each genomic
705   cytosine position in the appropriate sequence context within each 1 Mb window, every 1 Mb.  *Populus* methylation
706   data[46] was for Tree 13 branch 1, from GEO (https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE132939
707   2020). Local correlations between methylation levels and gene count were determined using methods from
708   Niederhuth, et al. [43] to maximize relevance of the comparison. Thus, using Bedtools, the genome was divided into
709   100 kbp windows with 50 kbp overlaps. Methylation for each 100 kbp window was from averaging 100 bp window
710   methylation levels (as above). Genes per window were counted with Bedtools intersect, requiring at least 50% of
711   the gene span to be inside the window. Correlation between gene count and methylation level was calculated with
712   R cor()'s Pearson method with incomplete observations dropped.

713   **Data availability**

714   Data are available at NCBI (GCA_001633185.3, additional accessions TBD), European Variation Archive
715   (accession TBD), the project website (valleyoak.ucla.edu), and the project genome browser
716   (genomes.mcdb.ucla.edu/cgi-bin/hgTracks?dg=queLob3).

717   **References**

718   1   Kremer, A. & Hipp, A. L. Oaks: an evolutionary success story. *New Phytologist* **226**, 987-1011 (2020).
719   2   Hipp, A. L., Manos, P. S. & Cavender-Bares, J. How oak trees evolved to rule the forests of the Northern
720       Hemisphere. *Scientific American* **323**, 42-49, doi:10.1038/scientificamerican0820-42 (2020).
721   3   Barrón, E. *et al.* in *Oaks Physiological Ecology. Exploring the Functional Diversity of Genus Quercus L.*  (ed
722       Peguero-Pina J. Gil-Pelegrín E., Sancho-Knapik, D.)  39-105 (Springer, 2017).
723   4   Denk, T., Grimm, G. W., Manos, P. S., Deng, M. & Hipp, A. L. in *Oaks Physiological Ecology. Exploring the
724       Functional Diversity of Genus Quercus L..* Vol. 7 *Tree Physiology* (eds E. Gil-Pelegrín, J. Peguero-Pina, & D.
725       Sancho-Knapik)  (Springer, 2017).
726   5   Cavender-Bares, J. Diversity, distribution, and ecosystem services of the North American oaks. *Journal of
727       International Oak Society* **27**, 37-48 (2016).

728    6    Plomion, C. *et al.* Oak genome reveals facets of long lifespan. *Nature Plants* **4**, 440-452, doi:10.1038/s41477-
729         018-0172-3 (2018).
730    7    Salojärvi, J. *et al.* Genome sequencing and population genomic analyses provide insights into the adaptive
731         landscape of silver birch. *Nature Genetics* **49**, 904-912, doi:10.1038/ng.3862 (2017).
732    8    Li, Q. *et al.* RNA-directed DNA methylation enforces boundaries between heterochromatin and euchromatin
733         in the maize genome. *Proceedings of the National Academy of Sciences*, 201514680,
734         doi:10.1073/pnas.1514680112 (2015).
735    9    Sork, V. L. *et al.* First draft assembly and annotation of the genome of a California endemic oak. *Quercus
736         lobata* Née (Fagaceae). *G3: Genes | Genomes | Genetics* **11**, 3485-3495 (2016).
737    10   Putnam, N. H. *et al.* Chromosome-scale shotgun assembly using an in vitro method for long-range linkage.
738         *Genome Research* **26**, 342-350, doi:10.1101/gr.193474.115 (2016).
739    11   Bodénès, C., Chancerel, E., Ehrenmann, F., Kremer, A. & Plomion, C. High-density linkage mapping and
740         distribution of segregation distortion regions in the oak genome. *DNA Research* **23**, 115-124,
741         doi:10.1093/dnares/dsw001 (2016).
742    12   Ramos, A. M. *et al.* The draft genome sequence of cork oak. *Scientific Data* **5**, 180069,
743         doi:10.1038/sdata.2018.69 (2018).
744    13   Hipp, A. L. *et al.* Genomic landscape of the global oak phylogeny. *New Phytologist* **226**, 1198-1212 (2020).
745    14   Tuskan, G. A. *et al.* The genome of black cottonwood, *Populus trichocarpa* (Torr. &amp; Gray). *Science* **313**,
746         1596-1604, doi:10.1126/science.1128691 (2006).
747    15   Schiffels, S. & Durbin, R. Inferring human population size and separation history from multiple genome
748         sequences. *Nature Genetics* **46**, 919-925, doi:10.1038/ng.3015 (2014).
749    16   Rundel, P. W. *et al.* Mediterranean Biomes: Evolution of Their Vegetation, Floras, and Climate. *Annual Review
750         of Ecology, Evolution, and Systematics* **47**, 383-407, doi:10.1146/annurev-ecolsys-121415-032330 (2016).
751    17   Corbett-Detig, R. B., Hartl, D. L. & Sackton, T. B. Natural selection constrains neutral diversity across a wide
752         range of species. *PLOS Biology* **13**, e1002112, doi:10.1371/journal.pbio.1002112 (2015).
753    18   Myburg, A. A. *et al.* The genome of *Eucalyptus grandis*. *Nature* **510**, 356-362, doi:10.1038/nature13308
754         (2014).
755    19   Argout, X. *et al.* The genome of *Theobroma cacao*. *Nature Genetics* **43**, 101-108, doi:10.1038/ng.736 (2011).
756    20   Denoeud, F. *et al.* The coffee genome provides insight into the convergent evolution of caffeine biosynthesis.
757         *Science* **345**, 1181-1184, doi:10.1126/science.1255274 (2014).
758    21   Griffin, J. R. & Critchfield, W. B. *The distribution of the forest trees in California*.  (Pacific SW Forest and Range
759         Experiment Station,  U.S. Department of Agriculture Forest Service, 1972).
760    22   Stanke, M., Schöffmann, O., Morgenstern, B. & Waack, S. Gene prediction in eukaryotes with a generalized
761         hidden Markov model that uses hints from external sources. *BMC Bioinformatics* **7**, 62, doi:10.1186/1471-
762         2105-7-62 (2006).
763    23   El-Gebali, S. *et al.* The Pfam protein families database in 2019. *Nucleic Acids Research* **47**, D427-D432,
764         doi:10.1093/nar/gky995 (2018).
765    24   Seppey, M., Manni, M. & Zdobnov, E. M. in *Gene Prediction. Methods in Molecular Biology* Vol. 1962  (ed M.
766         Kollmar)  ( Humana, 2019).
767    25   Gururani, M. A. *et al.* Plant disease resistance genes: current status and future directions. *Physiological and
768         molecular plant pathology* **78**, 51-65 (2012).
769    26   Manzanares, C. *et al.* A gene encoding a DUF247 domain protein cosegregates with the S Self-Incompatibility
770         locus in perennial ryegrass. *Molecular Biology and Evolution* **33**, 870-884, doi:10.1093/molbev/msv335
771         (2015).
772    27   Thorogood, D. *et al.* A novel multivariate approach to phenotyping and association mapping of multi-locus
773         gametophytic self-incompatibility reveals s, z, and other loci in a perennial ryegrass (Poaceae) population.
774         *Frontiers in Plant Science* **8**, doi:10.3389/fpls.2017.01331 (2017).
775    28   Fujii, S., Kubo, K.-i. & Takayama, S. Non-self- and self-recognition models in plant self-incompatibility. *Nature
776         Plants* **2**, 16130, doi:10.1038/nplants.2016.130 (2016).
777    29   Iwano, M. & Takayama, S. Self/non-self discrimination in angiosperm self-incompatibility. *Curr. Opin. Plant
778         Biol.* **15**, 78-83, doi:https://doi.org/10.1016/j.pbi.2011.09.003 (2012).
779    30   Kubo, K.-i. *et al.* Gene duplication and genetic exchange drive the evolution of S-RNase-based self-
780         incompatibility in Petunia. *Nature Plants* **1**, 14005, doi:10.1038/nplants.2014.5 (2015).

781  31  Li, W. & Chetelat, R. T. A pollen factor linking inter- and intraspecific pollen rejection in tomato. *Science* **330**,
782      1827-1830, doi:10.1126/science.1197908 (2010).
783  32  De Franceschi, P., Dondini, L. & Sanzol, J. Molecular bases and evolutionary dynamics of self-incompatibility in
784      the Pyrinae (Rosaceae). *Journal of Experimental Botany* **63**, 4015-4032, doi:10.1093/jxb/ers108 (2012).
785  33  Li, M. *et al.* Genome structure and evolution of *Antirrhinum majus* L. *Nature Plants* **5**, 174-183,
786      doi:10.1038/s41477-018-0349-9 (2019).
787  34  Aguiar, B. *et al.* Convergent evolution at the gametophytic self-incompatibility system in *Malus* and *Prunus*.
788      *PLOS ONE* **10**, e0126138, doi:10.1371/journal.pone.0126138 (2015).
789  35  Zhang, Z. *et al.* A PECTIN METHYLESTERASE gene at the maize Ga1 locus confers male function in unilateral
790      cross-incompatibility. *Nature Communications* **9**, 3678, doi:10.1038/s41467-018-06139-8 (2018).
791  36  Chen, M. *et al.* EbARC1, an E3 ubiquitin ligase gene in *Erigeron breviscapus*, confers self-incompatibility in
792      transgenic *Arabidopsis thaliana*. *International Journal of Molecular Sciences* **21**, 1458 (2020).
793  37  Mistry, J. *et al.* Pfam: The protein families database in 2021. *Nucleic Acids Research* **49**, D412-D419,
794      doi:10.1093/nar/gkaa913 (2020).
795  38  Chanderbali, A. S., Berger, B. A., Howarth, D. G., Soltis, D. E. & Soltis, P. S. Evolution of floral diversity:
796      genomics, genes and gamma. *Philosophical Transactions of the Royal Society B-Biological Sciences* **372**,
797      doi:10.1098/rstb.2015.0509 (2017).
798  39  Jiao, Y. N. *et al.* A genome triplication associated with early diversification of the core eudicots. *Genome*
799      *Biology* **13**, doi:10.1186/gb-2012-13-1-r3 (2012).
800  40  Vekemans, D. *et al.* Gamma paleohexaploidy in the stem lineage of core eudicots: significance for MADS-Box
801      gene and species diversification. *Molecular Biology and Evolution* **29**, 3793-3806,
802      doi:10.1093/molbev/mss183 (2012).
803  41  Higo, A. *et al.* DNA methylation is reconfigured at the onset of reproduction in rice shoot apical meristem.
804      *Nature communications* **11**, 4079-4079, doi:10.1038/s41467-020-17963-2 (2020).
805  42  Gent, J. I. *et al.* CHH islands: de novo DNA methylation in near-gene chromatin regulation in maize. *Genome*
806      *Research* **23**, 628-637, doi:10.1101/gr.146985.112 (2013).
807  43  Niederhuth, C. E. *et al.* Widespread natural variation of DNA methylation within angiosperms. *Genome*
808      *biology* **17**, 194 (2016).
809  44  Gouil, Q. & Baulcombe, D. C. DNA methylation signatures of the plant chromomethyltransferases. *PLOS*
810      *Genetics* **12**, e1006526, doi:10.1371/journal.pgen.1006526 (2016).
811  45  Song, X. & Cao, X. Context and complexity: Analyzing methylation in trinucleotide sequences. *Trends in Plant*
812      *Science* **22**, 351-353, doi:https://doi.org/10.1016/j.tplants.2017.03.013 (2017).
813  46  Hofmeister, B. T. *et al.* A genome assembly and the somatic genetic and epigenetic mutation rate in a wild
814      long-lived perennial *Populus trichocarpa*. *Genome Biology* **21**, 259, doi:10.1186/s13059-020-02162-5 (2020).
815  47  Carpentier, M.-C. *et al.* Retrotranspositional landscape of Asian rice revealed by 3000 genomes. *Nature*
816      *Communications* **10**, 24, doi:10.1038/s41467-018-07974-5 (2019).
817  48  Choi, J. Y. & Lee, Y. C. G. Double-edged sword: The evolutionary consequences of the epigenetic silencing of
818      transposable elements. *PLOS Genetics* **16**, e1008872, doi:10.1371/journal.pgen.1008872 (2020).
819  49  Mei, W., Stetter, M. G., Gates, D. J., Stitzer, M. C. & Ross-Ibarra, J. Adaptation in plant genomes: Bigger is
820      different. *Am J Bot* **105**, 16-19, doi:https://doi.org/10.1002/ajb2.1002 (2018).
821  50  Crisp, P. A. *et al.* Stable unmethylated DNA demarcates expressed genes and their cis-regulatory space in
822      plant genomes. *Proceedings of the National Academy of Sciences of the United States of America* **117**, 23991-
823      24000, doi:10.1073/pnas.2010250117 (2020).
824  51  Ricci, W. A. *et al.* Widespread long-range cis-regulatory elements in the maize genome. *Nature Plants* **5**, 1237-
825      1249, doi:10.1038/s41477-019-0547-0 (2019).
826  52  Rodgers-Melnick, E., Vera, D. L., Bass, H. W. & Buckler, E. S. Open chromatin reveals the functional maize
827      genome. *Proceedings of the National Academy of Sciences*, 201525244, doi:10.1073/pnas.1525244113
828      (2016).
829  53  Baduel, P., Quadrana, L., Hunter, B., Bomblies, K. & Colot, V. Relaxed purifying selection in autopolyploids
830      drives transposable element over-accumulation which provides variants for local adaptation. *Nature*
831      *Communications* **10**, 5818, doi:10.1038/s41467-019-13730-0 (2019).
832  54  Chuong, E. B., Elde, N. C. & Feschotte, C. Regulatory activities of transposable elements: from conflicts to
833      benefits. *Nature Reviews Genetics* **18**, 71-86, doi:10.1038/nrg.2016.139 (2017).

834 55 Li, E. *et al.* Long-range interactions between proximal and distal regulatory regions in maize. *Nature*
835 *Communications* **10**, 2633, doi:10.1038/s41467-019-10603-4 (2019).
836 56 Vanrobays, E., Thomas, M., Tatout, C. in *Annual Plant Reviews online* (ed J.A. Roberts) 157-190 (2017).
837 57 Makarevitch, I. *et al.* Genomic distribution of maize facultative heterochromatin marked by trimethylation of
838 H3K27. *The Plant Cell* **25**, 780-793, doi:10.1105/tpc.112.106427 (2013).
839 58 Zhao, L. *et al.* Chromatin loops associated with active genes and heterochromatin shape rice genome
840 architecture for transcriptional regulation. *Nature Communications* **10**, 3640, doi:10.1038/s41467-019-11535-
841 9 (2019).
842 59 Long, J. C. *et al.* Decrease in DNA methylation 1 (DDM1) is required for the formation of mCHH islands in
843 maize. *Journal of Integrative Plant Biology* **61**, 749-764, doi:https://doi.org/10.1111/jipb.12733 (2019).
844 60 Achour, Z. *et al.* Low temperature triggers genome-wide hypermethylation of transposable elements and
845 centromeres in maize. *bioRxiv*, 573915, doi:10.1101/573915 (2019).
846 61 Sasaki, E., Kawakatsu, T., Ecker, J. R. & Nordborg, M. Common alleles of CMT2 and NRPE1 are major
847 determinants of CHH methylation variation in *Arabidopsis thaliana*. *PLOS Genetics* **15**, e1008492,
848 doi:10.1371/journal.pgen.1008492 (2020).
849 62 Kear, P. J. & McClure, B. in *Adv Exp Med Biol* (ed Carlos López-Larrea) 108-123 (Springer US, 2012).
850 63 Boavida, L. C., Silva, J. P. & Feijo, J. A. Sexual reproduction in the cork oak (*Quercus sober* L). - II. Crossing
851 intra- and interspecific barriers. *Sex. Plant Reprod.* **14**, 143-152, doi:10.1007/s004970100100 (2001).
852 64 Charlesworth, D., Vekemans, X., Castric, V. & Glémin, S. Plant self-incompatibility systems: a molecular
853 evolutionary perspective. *New Phytologist* **168**, 61-69, doi:https://doi.org/10.1111/j.1469-8137.2005.01443.x
854 (2005).
855 65 Johnson, D. J., Beaulieu, W. T., Bever, J. D. & Clay, K. Conspecific negative density dependence and forest
856 diversity. *Science* **336**, 904-907, doi:10.1126/science.1220269 (2012).
857 66 Bever, J. D., Mangan, S. A. & Alexander, H. M. Maintenance of plant species diversity by pathogens. *Annual*
858 *Review of Ecology, Evolution, and Systematics* **46**, 305-325, doi:10.1146/annurev-ecolsys-112414-054306
859 (2015).
860 67 Marden, J. H. *et al.* Ecological genomics of tropical trees: how local population size and allelic diversity of
861 resistance genes relate to immune responses, cosusceptibility to pathogens, and negative density
862 dependence. *Molecular Ecology* **26**, 2498-2513, doi:https://doi.org/10.1111/mec.13999 (2017).
863 68 Stump, S. M., Marden, J. H., Beckman, N. G., Mangan, S. A. & Comita, L. S. Resistance genes affect how
864 pathogens maintain plant abundance and diversity. *The American Naturalist* **196**, 472-486,
865 doi:10.1086/710486 (2020).
866 69 Xue, J.-Y., Takken, F. L. W., Nepal, M. P., Maekawa, T. & Shao, Z.-Q. Editorial: Evolution and Functional
867 Mechanisms of Plant Disease Resistance. *Frontiers in Genetics* **11**, doi:10.3389/fgene.2020.593240 (2020).
868 70 Karasov, T. L., Shirsekar, G., Schwab, R. & Weigel, D. What natural variation can teach us about resistance
869 durability. *Curr. Opin. Plant Biol.* **56**, 89-98, doi:https://doi.org/10.1016/j.pbi.2020.04.010 (2020).
870 71 Rensing, S. A. Gene duplication as a driver of plant morphogenetic evolution. *Curr. Opin. Plant Biol.* **17**, 43-48,
871 doi:10.1016/j.pbi.2013.11.002 (2014).
872 72 Defoort, J., Van de Peer, Y. & Carretero-Paulet, L. The evolution of gene duplicates in angiosperms and the
873 impact of protein–protein interactions and the mechanism of duplication. *Genome Biology and Evolution* **11**,
874 2292-2305, doi:10.1093/gbe/evz156 (2019).
875 73 Sork, V., Dyer, R., Davis, F. & Smouse., P. in *Proceedings of the Fifth Symposium on Oak Woodlands: Oaks in*
876 *California's Changing Landscape. 2001 October 22-25; San Diego, CA. Gen. Tech. Rep. PSW-GTR-184.* (ed R.B.
877 McCreary Standiford, D; Purcell, K.L.) 427-444 (Pacific Southwest Research Station, Forest Service, U.S.
878 Department of Agriculture; , 2002).
879 74 Pluess, A. R. *et al.* Short distance pollen movement in a wind-pollinated tree, *Quercus lobata* (Fagaceae).
880 *Forest Ecology and Management* **258**, 735-744, doi:10.1016/j.foreco.2009.05.014 (2009).
881 75 Sork, V. L. & Smouse, P. E. Genetic analysis of landscape connectivity in tree populations. *Landscape Ecology*
882 **21**, 821-836 (2006).
883 76 Sork, V. L., Smouse, P. E., Grivet, D. & Scofield, D. G. Impact of asymmetric male and female gamete dispersal
884 on allelic diversity and spatial genetic structure in valley oak (*Quercus lobata* Née). *Evolutionary ecology* **29**,
885 927-945 (2015).

886  77  Grivet, D., Deguilloux, M.-F., Petit, R. J. & Sork, V. L. Contrasting patterns of historical colonization in white
887      oaks (*Quercus* spp.) in California and Europe. *Molecular Ecology* **15**, 4085-4093 (2006).
888  78  Gugger, P. F., Ikegami, M. & Sork, V. L. Influence of late Quaternary climate change on present patterns of
889      genetic variation in valley oak, *Quercus lobata* Née. *Molecular Ecology* **22**, 3598-3612,
890      doi:10.1111/mec.12317 (2013).
891  79  Petit, R. J. *et al.* Chloroplast DNA footprints of postglacial recolonization by oaks. *Proceedings of the National*
892      *Academy of Sciences of the United States of America* **94**, 9996-10001 (1997).
893  80  Anderson, M. K. *Tending the Wild: Native American Knowledge and the Management of California's Natural*
894      *Resources*.  555 (University of California Press, 2005).
895  81  Whipple, A. A., Grossinger, R. M. & Davis, F. W. Shifting baselines in a California oak savanna: Nineteenth
896      century data to inform restoration scenarios. *Restoration Ecology* **19**, 88-101, doi:10.1111/j.1526-
897      100X.2009.00633.x (2011).
898  82  McLaughlin, B. C. & Zavaleta, E. S. Predicting species responses to climate change: demography and climate
899      microrefugia in California valley oak (*Quercus lobata*). *Global Change Biology* **18**, 2301-2312,
900      doi:10.1111/J.1365-2486.2011.02630.X (2012).
901  83  Tyler, C. M., Kuhn, B. & Davis, F. W. Demography and recruitment limitations of three oak species in
902      California. *Quarterly Review of Biology* **81**, 127-152 (2006).
903  84  Sork, V. L. *et al.* Gene movement and genetic association with regional climate gradients in California valley
904      oak (*Quercus lobata* Née) in the face of climate change. *Molecular Ecology* **19**, 3806-3823,
905      doi:10.1111/j.1365-294X.2010.04726.x (2010).
906  85  Kueppers, L. N., Snyder, M. A., Sloan, L. C., Zavaleta, E. S. & Fulfrost, B. Modeled regional climate change and
907      California endemic oak ranges. *Proceedings of the National Academy of Sciences of the United States of*
908      *America* **102**, 16281-16286 (2005).
909  86  Zimin, A. V. *et al.* Hybrid assembly of the large and highly repetitive genome of *Aegilops tauschii*, a progenitor
910      of bread wheat, with the MaSuRCA mega-reads algorithm. *Genome Research* **27**, 787-792,
911      doi:10.1101/gr.213405.116 (2017).
912  87  Cokus, S. J., Gugger, P. F. & Sork, V. L. Evolutionary insights from *de novo* transcriptome assembly and SNP
913      discovery in California white oaks. *BMC Genomics* **16**, 552 (2015).
914  88  Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics*
915      **25**, 1754-1760, doi:10.1093/bioinformatics/btp324 (2009).
916  89  Lepoittevin, C. *et al.* Single-nucleotide polymorphism discovery and validation in high-density SNP array for
917      genetic analysis in European white oaks. *Molecular Ecology Resources* **15**, 1446-1459, doi:10.1111/1755-
918      0998.12407 (2015).
919  90  Kelleher, J., Etheridge, A. M. & McVean, G. Efficient Coalescent Simulation and Genealogical Analysis for Large
920      Sample Sizes. *PLOS Computational Biology* **12**, e1004842, doi:10.1371/journal.pcbi.1004842 (2016).
921  91  Beichman, A. C., Phung, T. N. & Lohmueller, K. E. Comparison of Single Genome and Allele Frequency Data
922      Reveals Discordant Demographic Histories. *G3: Genes|Genomes|Genetics* **7**, 3605-3620,
923      doi:10.1534/g3.117.300259 (2017).
924  92  Schultz, M. D. *et al.* Human body epigenome maps reveal noncanonical DNA methylation variation. *Nature*
925      **523**, 212-216 (2015).
926  93  Ramírez, F. *et al.* deepTools2: a next generation web server for deep-sequencing data analysis. *Nucleic Acids*
927      *Research* **44**, W160-W165, doi:10.1093/nar/gkw257 (2016).
928  94  Wickham, H. *ggplot2: Elegant Graphics for Data Analysis*.  (Springer-Velgag, 2016).

929  ### *Acknowledgments*

937  Expression Analysis Cores of the UC Davis Genome Center, supported by NIH Shared Instrumentation Grant
938  1S10OD010786-01; and whole genome resequencing, WGBS, and RNA-Seq for demographics, methylation, and
939  transcriptomic studies, respectively, were done at the UCLA Broad Stem Cell Genome Core facility. We thank Lily
940  Shiue and Thomas Swale, Dovetail Genomics, for the high-quality scaffolding performed by HiRise.

941  ### *Author contributions*

942  VLS, MP, and SLS conceived the overall project design and management and obtained grant support; VLS initiated,
943  coordinated, and supervised the project and manuscript. SJC annotated and analyzed genes and repeats, designed
944  and conducted genome comparative analyses, and created/wrote figures, results, methods, and supplementary
945  information (SI) these sections. STF-G analyzed methylomes, designed and conducted comparative methylation
946  analyses, created/wrote figures, results, discussion, methods, and supplementary information (SI) for this topic;
947  called genetic variants (GATK) for the demographic analysis; submitted genomic resources for public availability.
948  AVZ and DP assembled and curated the genome sequence and contributed to results, methods, and
949  supplementary information (SI) for these sections. JG and YZ analyzed genetic variation data. JG conducted
950  demographic analysis. KEL designed and supervised the demographic analysis and JG and KEL contributed results,
951  methods, and SI for this section. STF-G and SJC examined DUF247. CLH conducted lab preparation for DNA
952  sequencing, resequencing, bisulfite sequencing and RNA sequencing. VLS, SJC, STF-G, AVZ, JG, PFG, KEL, MP, and
953  SLS edited text. SJC edited manuscript figures. SJC, STF-G, and VLS interpreted data and wrote the manuscript.

954  ### *Competing interests*

955  The authors declare no competing interests.

956  ### *Additional information*

957  **Supplemental Information** (pdf)

958  **Auxiliary documents** (Auxiliary Spreadsheets 1 and 2, Excel document)

959