








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

Victoria L. Sork <sup>1,2,9</sup>✉, Shawn J. Cokus<sup>3,9</sup>, Sorel T. Fitz-Gibbon <sup>1,9</sup>, Aleksey V. Zimin<sup>4,5</sup>, Daniela Puiu<sup>4</sup>, Jesse A. Garcia<sup>1</sup>, Paul F. Guggenberger <sup>6</sup>, Claudia L. Henriquez<sup>1</sup>, Ying Zhen <sup>1</sup>, Kirk E. Lohmueller <sup>1,7</sup>, Matteo Pellegrini <sup>3</sup> & Steven L. Salzberg <sup>4,8</sup>

The genus *Quercus*, which emerged ~55 million years ago during globally warm temperatures, diversified into ~450 extant 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 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  $\gamma$  whole genome hexaploid duplication event and by tandem duplication within families, including numerous resistance genes and a very large block of duplicated DUF247 genes, which have been found to be associated with self-incompatibility in grasses. 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. Collectively, these features promote genetic and phenotypic variation that would facilitate adaptability to changing environments.

<sup>1</sup>Department of Ecology and Evolutionary Biology, University of California, Los Angeles, CA 90095-1438, USA. <sup>2</sup>Institute of the Environment and Sustainability, University of California, Los Angeles, CA 90095, USA. <sup>3</sup>Department of Molecular, Cell, and Developmental Biology, University of California, Los Angeles, CA 90095-7239, USA. <sup>4</sup>Center for Computational Biology, Whiting School of Engineering, Johns Hopkins University, Baltimore, MD 21218, USA. <sup>5</sup>Department of Biomedical Engineering, Johns Hopkins University, Baltimore, MD 21218, USA. <sup>6</sup>Appalachian Laboratory, University of Maryland Center for Environmental Science, Frostburg, MD 21532, USA. <sup>7</sup>Department of Human Genetics, David Geffen School of Medicine, University of California, Los Angeles, CA 90095, USA. <sup>8</sup>Departments of Biomedical Engineering, Computer Science, and Biostatistics, Johns Hopkins University, Baltimore, MD 21218, USA. <sup>9</sup>These authors contributed equally: Victoria L. Sork, Shawn J. Cokus, Sorel T. Fitz-Gibbon. ✉email: [vsork@ucla.edu](mailto:vsork@ucla.edu)

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)<sup>1,2</sup>. The genus evolved in the palearctic during a time when the earth experienced a warmer climate<sup>3</sup>. 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<sup>3</sup>. 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<sup>3</sup>. This genus, which has diversified into two subgenera, eight sections, and >450 species<sup>4</sup>, is an “evolutionary success story”<sup>1</sup>. In North America, oaks have more biomass than any other woody plant genus, including pines<sup>5</sup>, 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 (*Quercus 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 respond 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.<sup>6</sup> 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-recognition system that could prevent self-fertilization and promote outcrossing, or selectively allow occasional hybridizations. We also find many genes retained from the ancient  $\gamma$  paleohexaploid duplication event of the core eudicots. These are enriched for transcription factors and housekeeping genes, which may be more subject to strong (hard) selective sweeps than the tandemly duplicated genes<sup>7</sup>. Finally, we find some surprising similarities with the genomes of Poaceae (grasses—also highly successful plants). DNA methylation (BS-Seq) patterns indicate heterochromatin-rich chromosome arms and additionally show CHH methylation peaks upstream of the transcription start sites. Such prominent “mCHH islands” are known in maize<sup>8</sup> and a few other plants. These features could both affect gene expression and also facilitate tandem duplication events creating phenotypic variation and opportunities for selection.

## Results

**Genome assembly.** An initial draft genome (version 1.0)<sup>9</sup> was assembled from small ( $\approx 150\times$  coverage) and large insert ( $\approx 50\times$ ) Illumina paired-end reads. The final assembly (version 3.0) was constructed with the addition of Pacific Biosciences long reads ( $\approx 80\times$ ) and Hi-C long-range links produced by Dovetail Genomics and the HiRise re-scaffolder<sup>10</sup>, dramatically increasing NG50 scaffold size from 2 kbp to 75 Mbp (see “Methods”). The 12 longest scaffolds (chromosomes) were named and oriented to agree with pedunculate oak *Q. robur*<sup>6</sup>, and correspond (in order, but not generally orientation) with the 12 linkage groups (LGs) of an existing moderate-density linkage map of *Q. robur*  $\times$  *Q. petraea*<sup>11</sup>

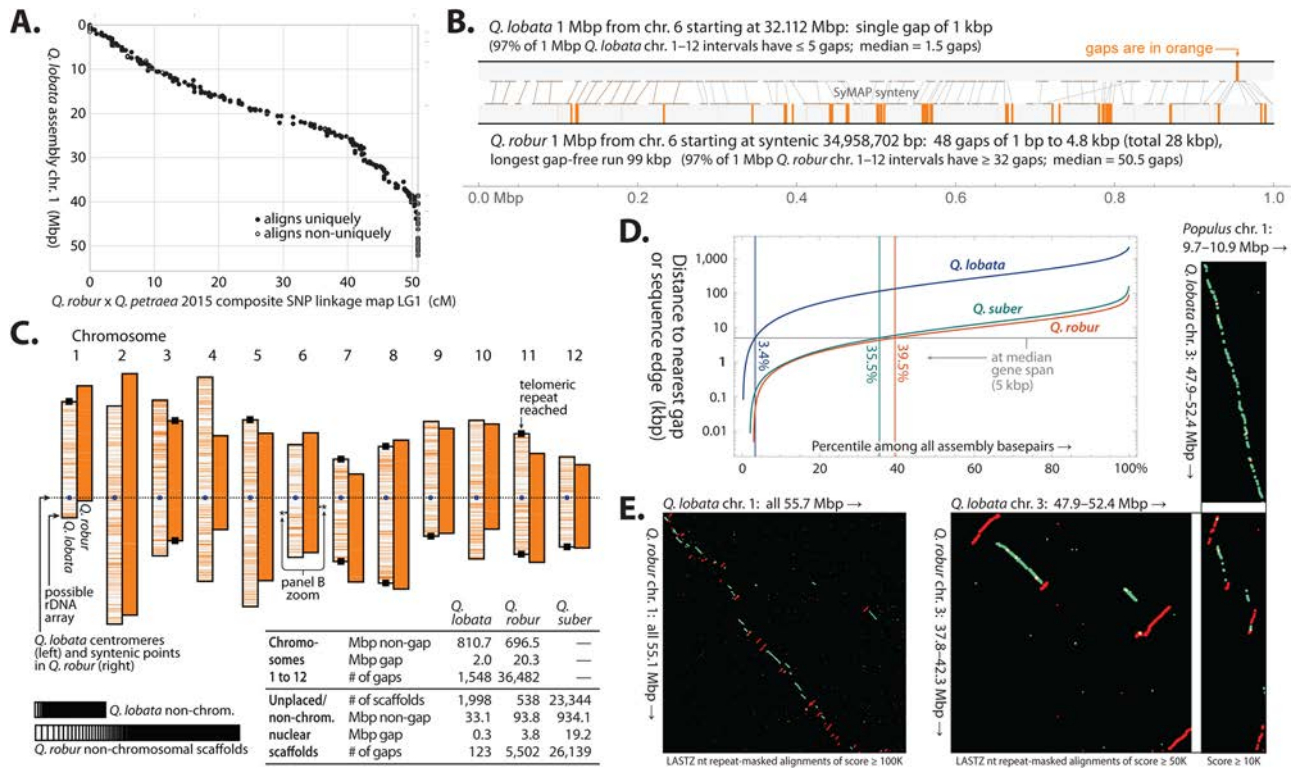
that we did not use during sequence construction. This linkage map consists of 4217 sequence context-defined single-nucleotide polymorphism (SNP) markers (after dropping 22 SNPs associated with 2 LGs each). The LGs and our chromosomes show a predominantly monotonic one-to-one correspondence (e.g., chr. 1 in Fig. 1A; see details in Supplementary Note 2. Validation and orientation of chromosomes, Supplementary Figs. 1 and 2, and Supplementary Table 2). We found that 99% of SNPs had at least one BLASTN alignment to our genome, and 98% of these had at least one alignment to the same chromosome as its LG. In addition, 95% of SNPs had all alignments to the same chromosome, and 86% had a unique alignment. A small stretch of our chromosome 1 was found to be a mis-assembled mitochondrial sequence and was replaced by a gap of the same length (Supplementary Fig. 3).

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

Our assembly is characterized by much higher contiguity than the other two. Comparing a representative 1 Mbp from *Q. lobata* and the syntenic 1 Mbp from *Q. robur* (Fig. 1B), the former has a single gap of 1 kbp while the latter has 28 kbp in 48 gaps of 1 bp to 5 kbp each. This pattern is typical: over all 1 Mbp regions from chr. 1–12, *Q. lobata* has median 1–2 gaps (97% of regions have  $\leq 5$ ), but *Q. robur* has 50–51 (97% having  $\geq 32$ ). Our assembly reaches telomeric repeats on both ends of four chromosomes and on one end of four more (Telomeric repeats, centromeres, and rDNA are discussed in Supplementary Note 6. “Repetitive sequences.”). Visualizing the entire *Q. lobata* and *Q. robur* assemblies (Fig. 1C), *Q. robur* gaps appear nearly solid. The percent of non-gap sequence placed in a chromosome is 96% in valley oak vs. 88% in pedunculate oak (and 0% in cork oak). More than a third of *Q. robur* and *Q. suber* base pairs are closer than a median gene span (5 kbp) to an assembly gap or sequence edge, while 96% of *Q. lobata* base pairs are further away (Fig. 1D).

Apparent segmental rearrangements and inversions between *Q. lobata* and *Q. robur* were unexpectedly prevalent (e.g., Fig. 1E left shows chr. 1 vs. chr. 1 as typical). Most of these, however, are likely scaffolding errors in *Q. robur*. Pedunculate oak has much smaller contigs, and its scaffolding was constructed using linkage maps (which are low in resolution compared to Hi-C) as well as synteny to *Prunus*, which may lead to mistakes in order and orientation of contigs (especially for small contigs). By contrast, alignments of *Q. lobata* with more distant species (*Populus*, *Eucalyptus*, *Theobroma*, and *Coffea*) showed numerous and widespread regions in continuous syntenies where *Q. robur* was not as continuous; to illustrate, Fig. 1E right shows Mbp-scale regions of chr. 3 of the two *Quercus* vs. chr. 1 of *Populus*. Further, a comparison of the formerly mentioned LGs from *Q. robur*  $\times$  *Q. petraea* to both the *Q. lobata* and *Q. robur* assemblies shows *Q. robur* with more disagreements (Supplementary Figs. 1 and 2). Thus, with the currently available *Q. robur* assembly, we conclude that alignments of *Q. robur* versus *Q. lobata* are not reflective of true rearrangements and inversions.

**Demographic histories of *Q. lobata* and *Q. robur*.** Ancient oaks evolved over 50 Mya, initially in the subtropical climate of the palearctic of the Northern Hemisphere and, as the planet cooled,



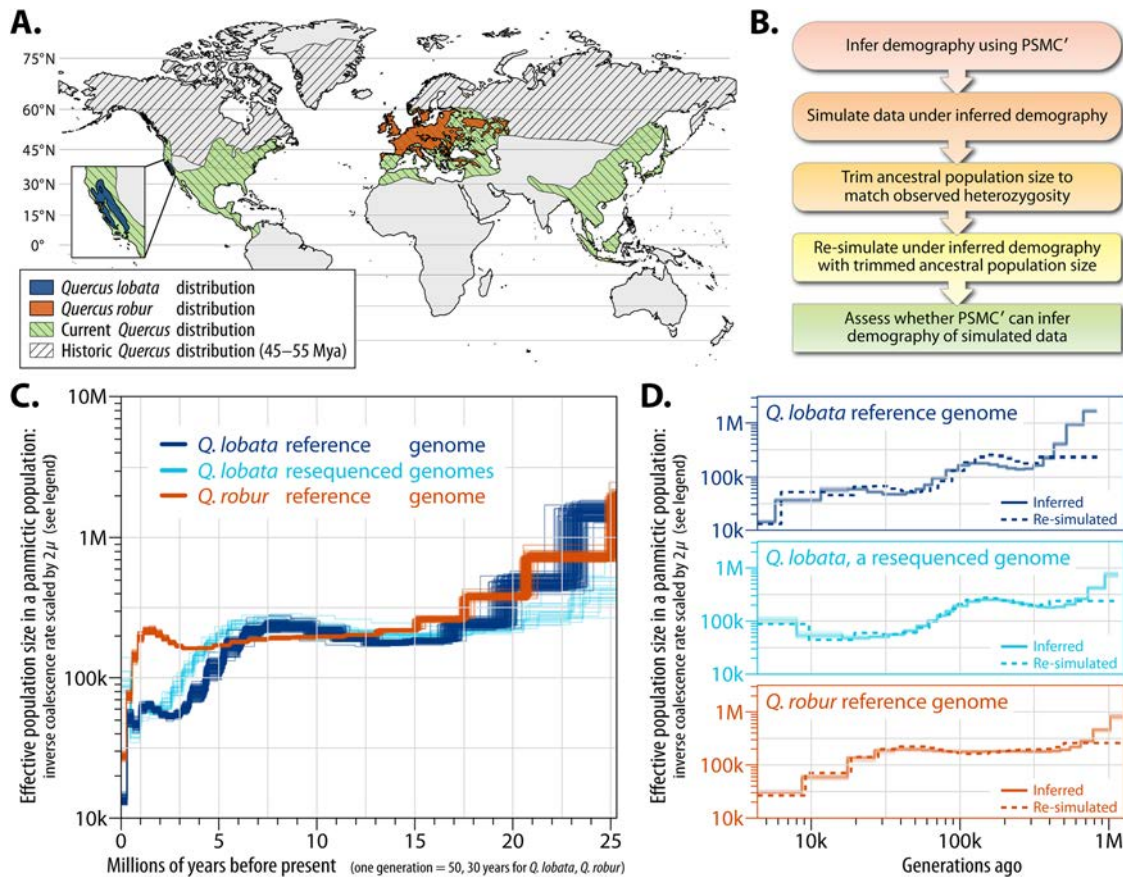
**Fig. 1 Overview of assemblies of *Q. lobata* tree SW786 (version 3.0), *Q. robur* (version PM1N)<sup>6</sup>, and *Q. suber* (version 1.0)<sup>12</sup>.** **A** Alignment of a linkage map linkage group 1 to *Q. lobata* chr. 1, exhibiting high concordance and overall monotonicity. **B** A representative 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*<sup>20</sup>, 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.

shifting southward to their contemporary distribution throughout the Northern Hemisphere (Fig. 2A). Consistent with the large range, we found heterozygosity (average 0.50%–0.66%; see Supplementary Note 3. Analysis of heterozygosity, and Supplementary Figs. 4 and 5) across the reference genome to be similar to but slightly less than the 0.73% computed for *Q. robur*, possibly due to the much larger species range of pedunculate oak and/or lower representation in the *Q. robur* assembly of highly homologous sequence loci resulting in an increased post-alignment pileup of multiple actual loci at single assembly loci. To gain insight into the population history of oaks, we inferred the effective population size ( $N_e$ ) of *Q. lobata* and *Q. robur* over time. The Pairwise Sequentially Markovian Coalescent (PSMC) method<sup>14</sup> applied to the individuals used to build the genomes mapped to their own assemblies (Fig. 2B) enabled examination of the last ≈25 M years of evolution (Supplementary Fig. 7). To verify accurate inference on this timespan, we generated simulated datasets using the inferred demographic history. We selected ancestral population sizes matching empirical heterozygosities in the reference genomes (Fig. 2C; see “Methods” and Supplementary Note 4. Demographic analysis, and Supplementary Figs. 8–10).

We ran PSMC on data simulated under trimmed demographic models and found the accurate inference of population size over time, except for the single oldest time step where population sizes were often over-estimated (Fig. 2D). The PSMC analysis indicates ancestral populations of both *Q. lobata* and *Q. robur* had high (>500k) effective population sizes that then showed

initially similar declines, perhaps as populations were shifting southward (Fig. 2C). *Q. lobata* shows an additional decline in  $N_e$  at ≈5 Mya, which coincides with a shift from a period of subtropical climate with year-round rainfall to a Mediterranean climate with summer drought<sup>15</sup>. By contrast, for *Q. robur* (being more widely distributed throughout Europe),  $N_e$  remained relatively flat until the last ≈1 M years. At this point, *Q. robur* declines to  $N_e < 50k$  (and below *Q. lobata*) during the “Ice Ages” when the region was experiencing a series of warm and cold periods creating genetic bottlenecks and expansions (Fig. 2C, D). Both species have retained sufficiently large effective population sizes to facilitate natural selection<sup>16</sup>. To determine if inference is affected by re-use of the same reads as used to build the reference assembly, we analyzed 19 re-sequenced *Q. lobata* individuals. These showed similar population size changes (Fig. 2C light blue) as SW786 (Fig. 2C dark blue), suggesting little bias from re-use.

**Repetitive sequences.** As with many plant species, the valley oak genome contains substantial repeats, with 54% of non-gap base pairs marked as repetitive by RepeatMasker in combination with a species-specific database constructed by RepeatModeler+Classifier (Fig. 3A, B; the modeling step was essential, as RepBase only marked 13%). The largest identified portion is transposable elements (TEs), primarily Copia and Gypsy elements of the long terminal repeat (LTR) type. The level of repetitiveness is similar to the 54% (disregarding gaps) found by application of the same process to *Q. robur* (for which Plomion et al.<sup>6</sup> reported 52% via REPET and other annotation, including manual curation).



**Fig. 2 Demographic evolutionary history analysis of *Q. lobata* and *Q. robur*.** **A** Historic and contemporary species ranges based on Barrón, Averyanova<sup>3</sup> and fossil occurrence records from the Global Biodiversity Information Facility website (<https://www.gbif.org/>, 19 January 2019). 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<sup>91</sup>. All data and map information are in the public domain. Land boundaries were obtained from NaturalEarthData.com (<https://www.naturalearthdata.com/about/terms-of-use/>). *Quercus* spp. records used to reconstruct the historic distribution were obtained from GBIF (GBIF.org, 19 January 2019, GBIF Occurrence Download). The current distribution of *Q. robur* was from EUFORGEN (<http://www.euforgen.org/species/quercus-robur/>). Occurrence records used to create the current distribution of *Quercus* and *Quercus lobata* were obtained from GBIF (GBIF.org, 22 January 2019, GBIF Occurrence Download). **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 Supplementary Note 4. Demographic analysis, and Supplementary Fig. 6 for other parameters). **D** PSMC' accurately infers demography (solid) on data simulated (dashed) under models fit to the empirical data.

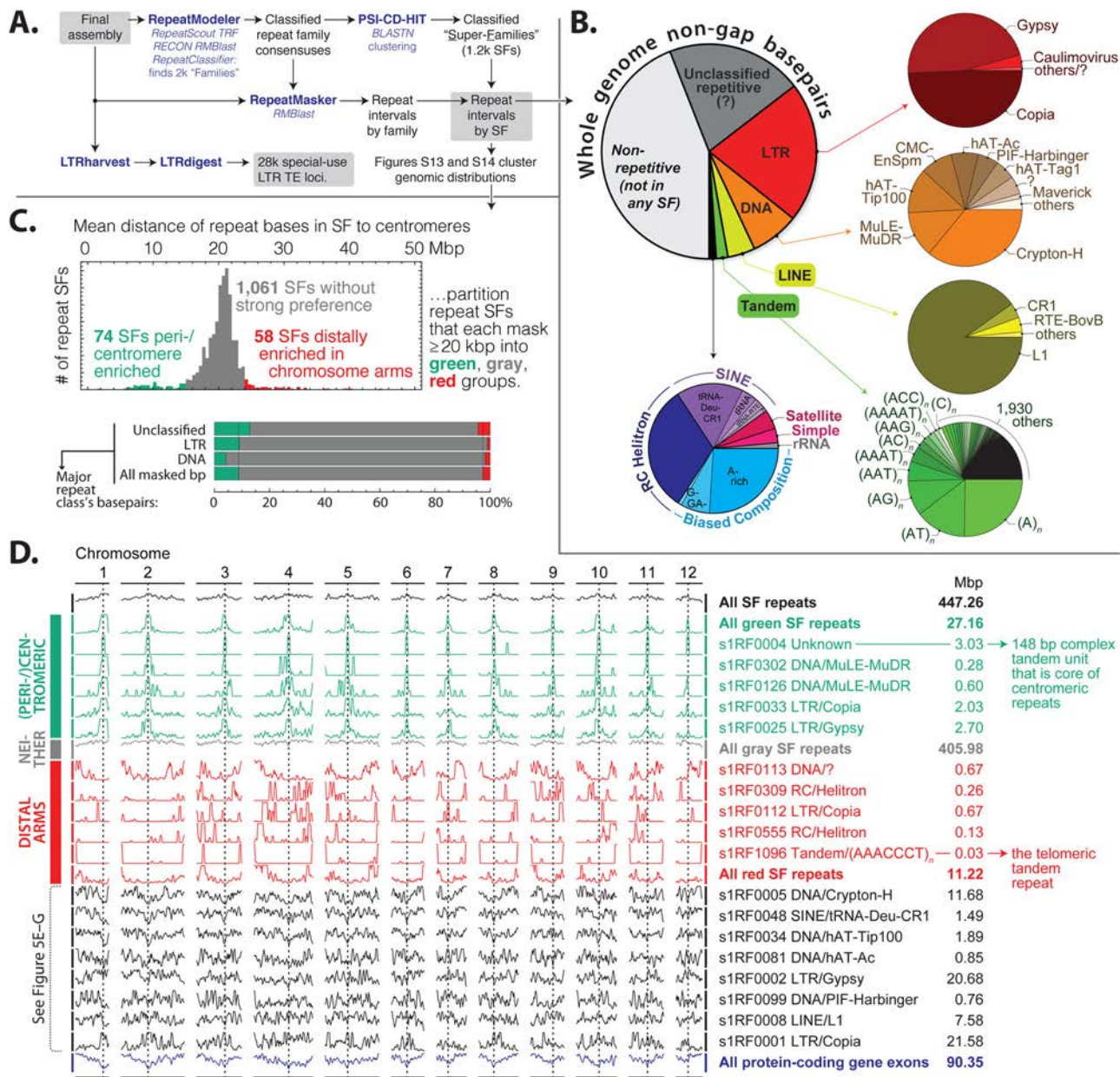
RepeatModeler+Classifier also detects 51% in *Q. suber*<sup>12</sup>, 55% in *Eucalyptus*<sup>17</sup>, 55% in *Theobroma*<sup>18</sup>, 51% in *Coffea*<sup>19</sup> and 43% in *Populus*<sup>20</sup>. Centromeric, telomeric, and rDNA repeats for valley oak were identified (see Supplementary Note 6. Repetitive sequences), and specific sequence-defined repeat superfamilies are correlated or anticorrelated to various levels with centromeric proximity, forming (as do protein-coding gene exons) density gradients that are the main chromosome-scale repeat-associated features, presumably reflecting overall chromatin structure (Supplementary Figs. 13 and 14 and Fig. 3C, D).

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.

**Gene prediction and annotation.** Using the AUGUSTUS gene modeler<sup>21</sup> 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 below “Methods” and Supplementary Note 7. Gene model statistics and possible R-genes in *Q. lobata*, *Q. robur*, and *Q. suber*, and Supplementary Table 3). As many corresponded to transposons with little expression or appeared hypothetical for other reasons, we removed 29k to obtain the primary set of 39,373 PCGs we report, of which 35k have at least one intron and all of which have UTRs annotated and are ostensibly complete. *Q. robur* reports only 29k PCG models, of which just 20k have introns, and about half UTRs; in the other direction, *Q. suber*'s annotation by NCBI (thinned to one isoform per locus) reports more 49k PCG loci (about half with UTRs), but a more comparable 36k with introns and 38k ostensibly complete, and with a much higher number containing transposon domains by comparison.

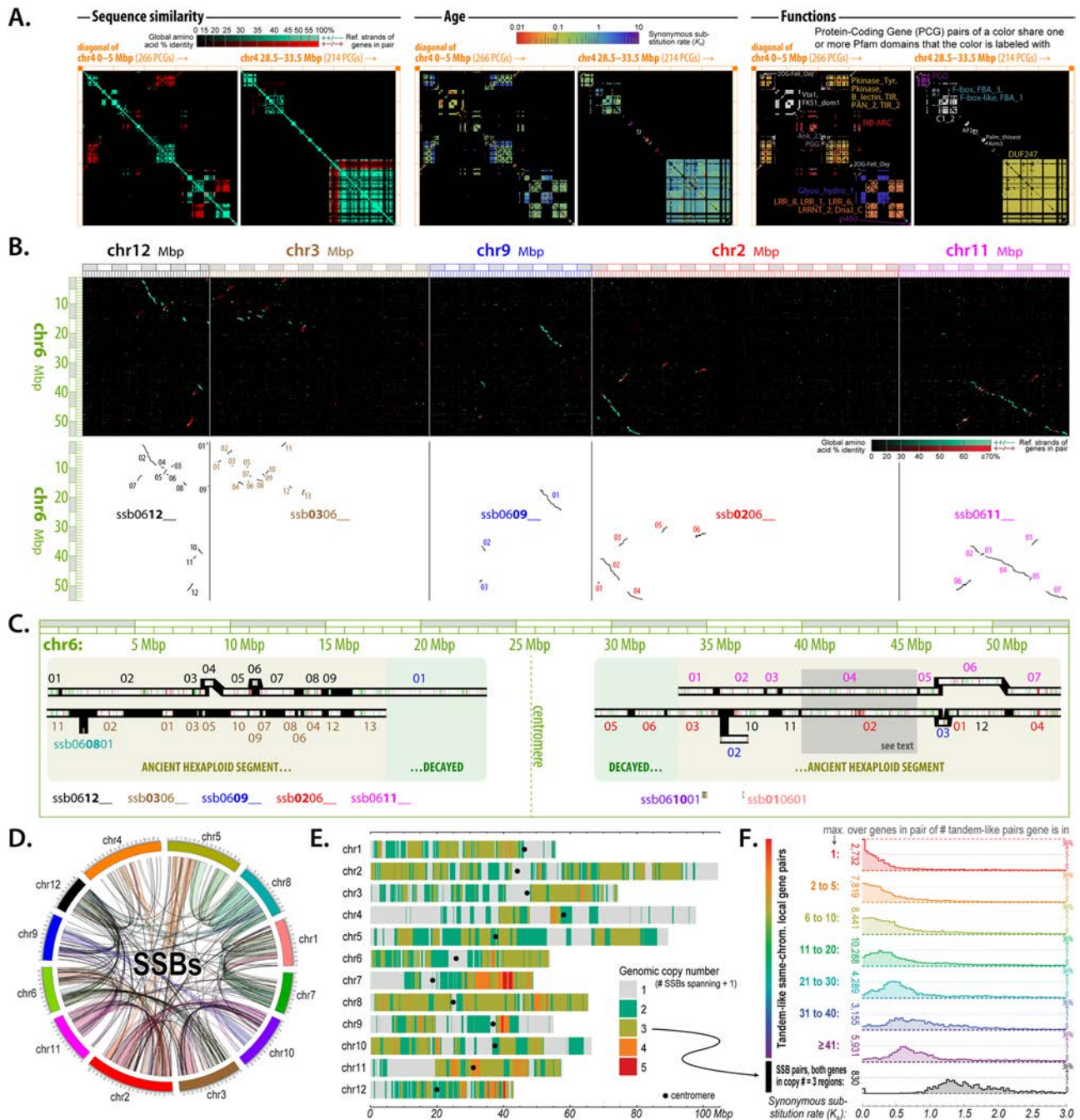
We assigned gene names, functions, and orthologs via the PANTHER and Pfam components of InterProScan, and OMA<sup>22</sup>. We evaluated the *Q. lobata*, *Q. robur*, and *Q. suber* scaffolds and single isoform PCG model sets with BUSCO (see below). *Q. lobata* compares favorably to the other two, and does not have the high multicopy anomaly of *Q. suber* in the 303-USCO ODB9 Eukaryota set<sup>23</sup>, or the high missing and fragmented fraction of *Q. robur*'s small model set (especially with the more comprehensive 2,121-USCO set for Eudicotyledons from ODB10).



**Fig. 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 1193 individual SFs each with ≥20 kbp distribute across chromosomes (Supplementary Figs. 13 and 14) 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, 1061 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 Fig. 5E-G.

**Gene duplications.** Protein–protein alignments among the *Q. lobata* PCGs exposed a rich panoply of duplication structure in terms of genomic positions, ages, and functions. Prominent and complex tandem-like blocks of high-similarity genes can be seen via visualizations of all-vs.-all alignments (see “Methods”). These duplications often involve local rearrangements, and can extend into megabases with dozens of genes involved at a time. Figure 4A (left third) exhibits two illustrative 5 Mbp regions of chr. 4. Approximately 40% of PCGs participate in these blocks, which have sizes of two to ≈100 genes each, with larger sizes rarified like a power law (Supplementary Fig. 15). Roughly a third of participating genes are duplicated only once, slightly more than half two to 20 times, and

only a tenth more than 20 times. Visualizations (e.g., coordinated Fig. 4A middle third) of the synonymous codon substitution rate ( $K_s$ ) over gene pairs in blocks suggest a wide variety of ages for the majority of retained expansion for individual blocks. Larger blocks tend to be older (Fig. 4F colored distributions), but even old blocks tend to have younger points suggestive of ongoing growth. While numerous tandem gene copies are shorter or have reduced or no RNA-Seq evidence of expression, many copies (even within larger blocks) are not particularly short or of lower expression and so do not appear to be pseudogenes. Functions of tandemly duplicated genes are diverse, as evident from the variety of Pfam domains they contain (e.g., coordinated Fig. 4A right third). Relatively few distinct



**Fig. 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. The largest block visible includes many DUF247 genes. 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  $\gamma$  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 ( $\geq \sim 1.0$ ) 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.

domains, however, are strongly enriched over all tandemly duplicated genes, and include NB-ARC, LRR\_8, B\_lectin, LRR\_1, TIR\_2, LRRNT\_2, p450, TIR, and PGG (associated with resistance/defense); Pkinase\_Tyr and Pkinase (signal transduction); UDPGT (the large UDP-glucuronosyl/glucosyl transferase family); S\_locus\_glycop,

PAN\_2, and DUF247 (see below); F-box, FBA\_3, and FBA\_1 (protein–protein interactions/degradation, signal transduction and regulation); and GST\_N, GST\_N\_3, and GST\_N\_2 (glutathione S-transferases, with functions including stress tolerance/signaling and detoxification).

Many of the strongly enriched domains are part of domain architectures of plant disease resistance genes (R-genes) that were identified by Gururani et al.<sup>24</sup>. It is difficult to be sure whether a putative R-gene is actually acting as a pathogen defense mechanism in a given plant species, but from the experimental evidence they reviewed, they identified eight classes of R-genes based on arrangements of domains and structural motifs. Using similar criteria for patterns (see Supplementary Note 7. Gene model statistics and possible R-genes in *Q. lobata*, *Q. robur*, and *Q. suber*), we analyzed the domain architecture of our 39,373 *Q. lobata* PCGs and identified 751 possible R-genes with strongly or highly suggestive patterns, and another 2176 with enriched potential (Supplementary Table 4). In the 25,808-gene *Q. robur* annotation<sup>6</sup>, we identify 632 + 1645 R-gene candidates, and for the 49,388-gene *Q. suber* set<sup>12</sup> (having been thinned to a single longest isoform per locus), we identify 723 + 2182 (Supplementary Table 4). (Because we examined Pfam domain hits to each project's respective gene models and did not, e.g., directly search the reference nucleotide sequences, differences here at least partially derive from variation in the annotation set quality/completeness.) Collectively, these amounts document likely high levels of R-genes in oak and indicate a tremendous opportunity for plant defense mechanisms.

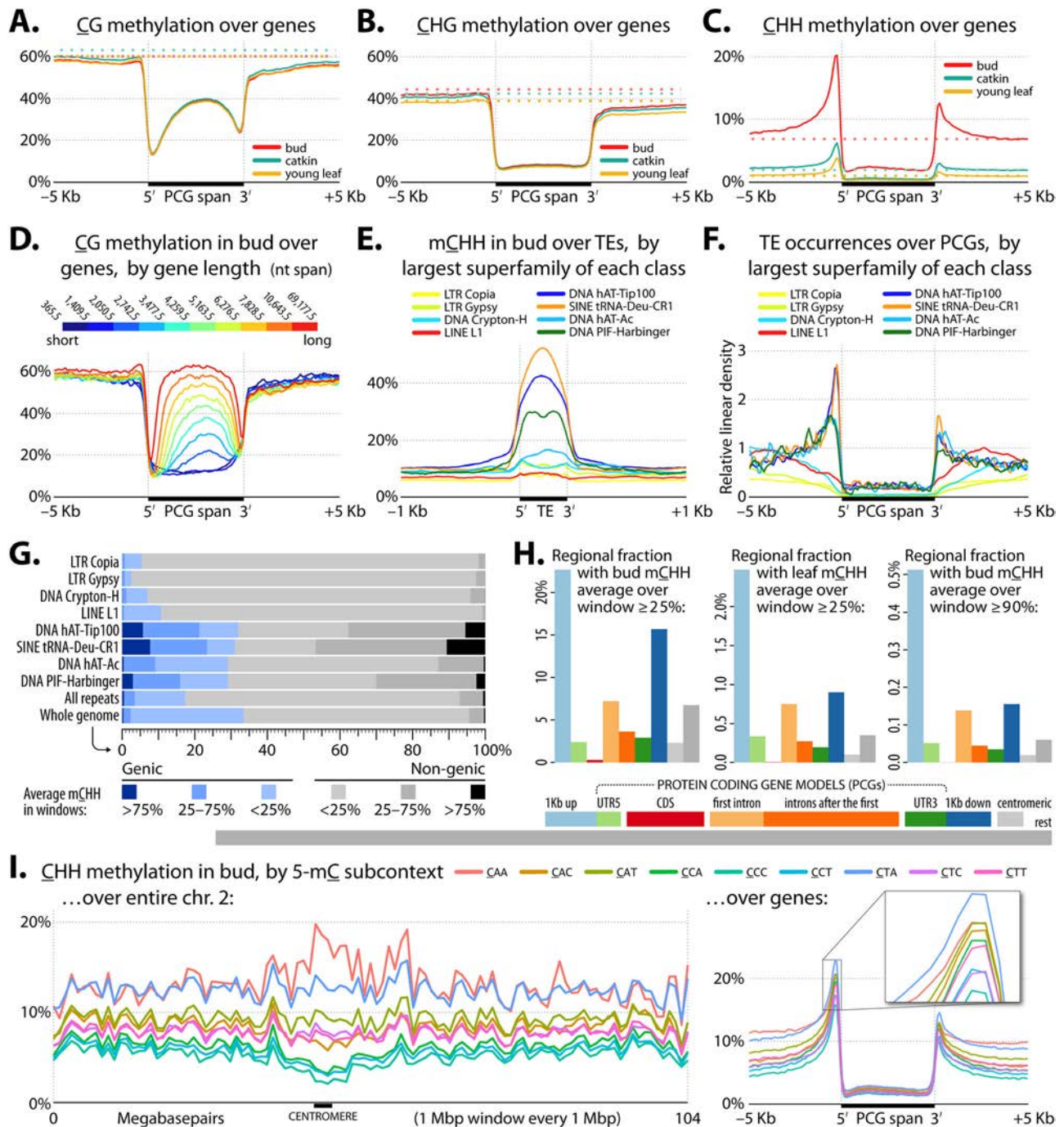
**Highly duplicated gene family, DUF247.** An investigation of PCGs found in blocks containing at least 30 tandemly duplicated genes uncovered DUF247 (PF03140) as the most enriched Pfam domain (see Supplementary Table 5). The span of the large DUF247-associated block (see the largest block in Fig. 4A) covers 34 predicted PCGs with a complete DUF247, 22 with partial DUF247's, and 17 additional genes. In total, we identified 186 DUF247-containing PCG matches (81 complete). DUF247 genes are entirely specific to plants and usually carry a single copy of the domain that comprises almost the entire coding sequence of the gene. Across the 104 plant genomes in Pfam Release 33.1<sup>25</sup> with  $\geq 17,500$  predicted protein entries (i.e., those most likely to be complete), the top five by the number of DUF247 occurrences are three tree species—*Juglans regia* (English walnut)  $n = 201$ ; *Eucalyptus grandis*,  $n = 188$ ; and *Populus trichocarpa* (black cottonwood),  $n = 161$ —and two polyploid cultivars (wheat,  $n = 192$ ; and peanut,  $n = 165$ ). One suggested function for DUF247 genes is reported in the Poaceae family, where two DUF247 genes in rye grass segregate with each of two known self-recognition loci, and are proposed to be the male determinants of a multi-locus self-incompatibility system<sup>26</sup>. Asparagus is another plant species where DUF247 is tied to a sex-related function<sup>27</sup>. If oak DUF247 genes are involved in reproductive recognition, a non-self-recognition system would be compatible with the high levels of duplication, as selection should favor the evolution of a complete set of proteins capable of recognizing the diverse array of S-haplotypes within the population<sup>28</sup>. Pairwise comparisons show variable and mostly high rates of divergence among DUF247 PCGs (see middle panel of Fig. 4A and Supplementary Fig. 11), and we observe general conservation of DUF247 blocks of genes between *Q. lobata* and *Q. robur*, another oak from the same section (Supplementary Fig. 12). These features suggest these genes have been selected for a range of stable functions. Further work will require careful separation of paralog/allelic diversity<sup>29</sup> and more cross-species comparisons, as well as experimental validations to reveal the functions of these prevalent and intriguing genes.

**Long-surviving duplicated genes.** Also striking in the visualizations of protein alignments were self-syntenic blocks (SSBs): syntenic runs of proteins within *Q. lobata*, generally between different chromosomes, with a variety of lengths and gene pair

densities. Figure 4B (top) shows chr. 6 vs. chr. 12/3/9/2/11 as exemplary (although in low resolution per limited space). For further analysis, 236 SSB runs, each with four to hundreds of gene pairs, were extracted (e.g., Fig. 4B bottom) and given accessions “ssbXXYYZZ” with  $XX \leq YY$  indicating the chromosomes involved and ZZ as a serial number; more than 7100 PCGs are directly involved. High-resolution examination made evident that, on any given chromosome, runs tended to end and begin close by, and for any particular point on a chromosome to be covered by very few runs (typically, zero to two), so that (nearly) disjoint SSBs could often be clearly ordered to form a small number of chromosome-scale chains (Fig. 4C black bars). While a few recent segmental duplications appear, most SSBs are likely long-lived syntelogs (i.e., homologous genes from the same ancient genome region) of the genome triplication polyploidy event  $\gamma$  associated with early diversification of the core eudicots, thought to have occurred about 120 Mya<sup>30–32</sup>. Several characteristics of the oak genome support the interpretation that these segmental duplications are derived from the ancient event  $\gamma$ . The high age of many SSBs is supported by the synonymous substitution rate ( $K_s$ ) for gene pairs in SSBs in triplicated regions being very high (almost entirely  $>1.0$ ; Fig. 4F black distribution). Moreover, the generally short length and scattered/scrambled nature of SSBs (which are within *Q. lobata*) compared to syntenies between *Q. lobata* and different species (*Populus*, *Eucalyptus*, *Theobroma*, and *Coffea*) would not be characteristic of recent events. Our direct findings are consistent with the *Q. robur* genome study<sup>6</sup>, which interpreted the  $K_s$  distribution of gene pairs for oak–peach orthologs vs. oak–oak to reveal duplicates resulting from the ancient  $\gamma$  triplication.

General triplication is clear from the gene pair-defined SSBs (e.g., Fig. 4C white bars, with green and red showing supporting gene pairs), but few syntenic gene triples have been retained. Thus, detection and characterization of the  $\gamma$  long-lived syntelogs would be unrepresentative for an analysis restricted to gene triples. For example, in the gray shaded region of Fig. 4C involving chr. 6/2/11 and spanning 320 chr. 6 genes, the 59 chr. 6 genes supporting local one-to-one chr. 6/2 syteny have only eleven chr. 6 genes in common with the 39 supporting local one-to-one chr. 6/11 syteny. Even before chaining as in Fig. 4C, two-thirds of the genome is in a SSB (Fig. 4D, E), with the largest fraction (34%) actually covered by two (consistent with triplication) and 27% by one (decayed triplications and a few recent segmental duplications); a third (34%) is not covered, and only 5% is covered by three or four SSBs (likely duplications post triplication). Relative to all genes, those in one or two gene pairs supporting SSBs tend to be of higher expression with a lower repetitive sequence in their immediate vicinity and are enriched for certain functional classes, including transcription factors and housekeeping genes (Supplementary Table 6), which is consistent with the gene balance hypothesis that genes duplicated during a polyploidization event are retained to allow retention of multiple interacting components of a complex that results in a multigenic phenotype<sup>33</sup>.

**Genome-wide patterns of DNA methylation and strong “mCHH islands”.** Whole-genome bisulfite sequencing for bud, catkin, and leaf tissue revealed mean DNA 5-methylcytosine methylation BS-Seq levels in CG (mCG) and CHG (mCHG) nucleotide contexts as relatively stable across tissues (Fig. 5A, B), while levels in CHH (mCHH; Fig. 5C) were notably higher in the bud than catkin and young leaf, likely due to the increased proportion of undifferentiated meristem tissue<sup>34</sup>. Mean levels for regions surrounding genes are similar to genome-wide means for all tissues in all contexts (mCHH 1–7%, mCHG 39–43%, mCG



**Fig. 5 DNA methylation in protein-coding genes and repeats.** **A–C** Average methylation levels (100 bp windows, introns included) 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. Supplementary Fig. 18 shows plots with introns excluded. **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 the 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): left side 1 Mbp windows across all of chr. 2, right side across genes (normalized to 5 kbp length) in bud tissue.

60–62%; Supplementary Fig. 16), with the exception of peaks of mCHH near transcription boundaries of genes (Fig. 5C). These mCHH peaks are similar in both position and scale above the background to the “mCHH islands” of maize<sup>35,36</sup> (Supplementary Data 1). We examined mCHH across representative repeat

superfamilies (SFs), specifically, those of highest mass within selected RepeatClassifier minor repeat classes, as seen in Fig. 5E. Within genic regions, three SFs—s1RF0048 (“SINE tRNA-Deu-CR1”), s1RF0034 (“DNA transposon hAT-Tip”), and s1RF0099 (“DNA transposon PIF-Harbinger”)—were both high in mCHH



and preferentially located in the highly methylated gene boundary regions (Fig. 5E, F). Members of these SFs are found in both genic and non-genic regions with broadly similar mCHH levels (Fig. 5G and Supplementary Fig. 18). However, in view of overall genome-wide mCHH levels (including centromeres and intergenic space), we find regions surrounding genes to be highly enriched for mCHH (Fig. 5H). Similar enrichment patterns are seen in bud and leaf, despite different overall mCHH levels (Fig. 5C, H), and similar patterns are also seen if mCHH window stringency is varied from 25% to 90%, although at these extremes we observe decreases in the relative amount of downstream and non-genic mCHH (Fig. 5H). All methylation is typically low near transcription boundaries (Fig. 5A), and remains low for mCHG and mCHH across gene bodies. However, gene body mCG rises for longer genes, reaching near-background levels in the longest genes (Fig. 5D).

**Broad distribution of heterochromatin.** *Q. lobata* appears to have heterochromatin dispersed throughout chromosomes more or less equally, with only a minor increase of density toward centromeres. This interpretation is based on both the distribution of genes and repeats as well as indications of widespread CMT-mediated DNA methylation, a pattern more similar to maize and rice methylomes than to the *Arabidopsis* and tomato methylomes in which the methylated repeats are concentrated in pericentromeric heterochromatic regions<sup>37,38</sup>. As such, a majority of repeat mass does not show a strong positional correlation with centromeres (Fig. 3D gray). Also, 92% of PCGs have a RepeatMasker-defined repeat within the gene's upstream 2 kbp, which is high, because, among 34 angiosperms, reported numbers range from 29% (*Arabidopsis*) to 94% (*Zea mays*), with an average of 50%<sup>36</sup>. (See also Supplementary Data 1.)

The second indication of heterochromatin-rich chromosome arms is the type of methylation found on intergenic repeats. Different mechanisms of generating plant mCHH, such as RNA-directed DNA Methylation (RdDM) or CMT2-mediated methylation, have been shown to have distinct preferences for specific nucleotide subcontexts (finer than CG/CHG/CHH: CAA vs. CAC vs. CAT, etc.). CMT-mediated mechanisms are typically responsible for methylation of heterochromatin and have much stronger biases than RdDM<sup>37</sup>. *Q. lobata* has strong CHH subcontext preferences at chromosome scale (Fig. 5I left and Supplementary Fig. 19). Bias patterns around centromeres (increased CAA and decreased CCH) are likely to indicate the general methylation pattern of heterochromatin in oaks, while chromosome arms represent a mix of genes and intergenic spaces. The peaks of mCHH surrounding gene boundaries (i.e., the possible mCHH islands) show another distinct pattern, with a preference for CAA strongly reduced (Fig. 5I right). Moving from gene boundaries toward intergenic space, the subcontext pattern progressively reverts to the likely heterochromatic signal of the centromeres (Supplementary Fig. 20).

An additional measure of the similarity in genome organization between oaks and grasses is the level of correlation between methylation and gene count across chromosomes. When we augment previously published Fig. 2D from Niederhuth et al.<sup>36</sup> with oak findings, oak is again found comparable to the grasses (Poaceae) and less to the other studied angiosperms (Fig. 6A). The low mCHH and gene count correlation reflects a combination of the unusually strong gene boundary mCHH relative to the background mCHH level (Supplementary Fig. 21 and Supplementary Data 1) and low average gene density in chromosome arms (Supplementary Fig. 22).

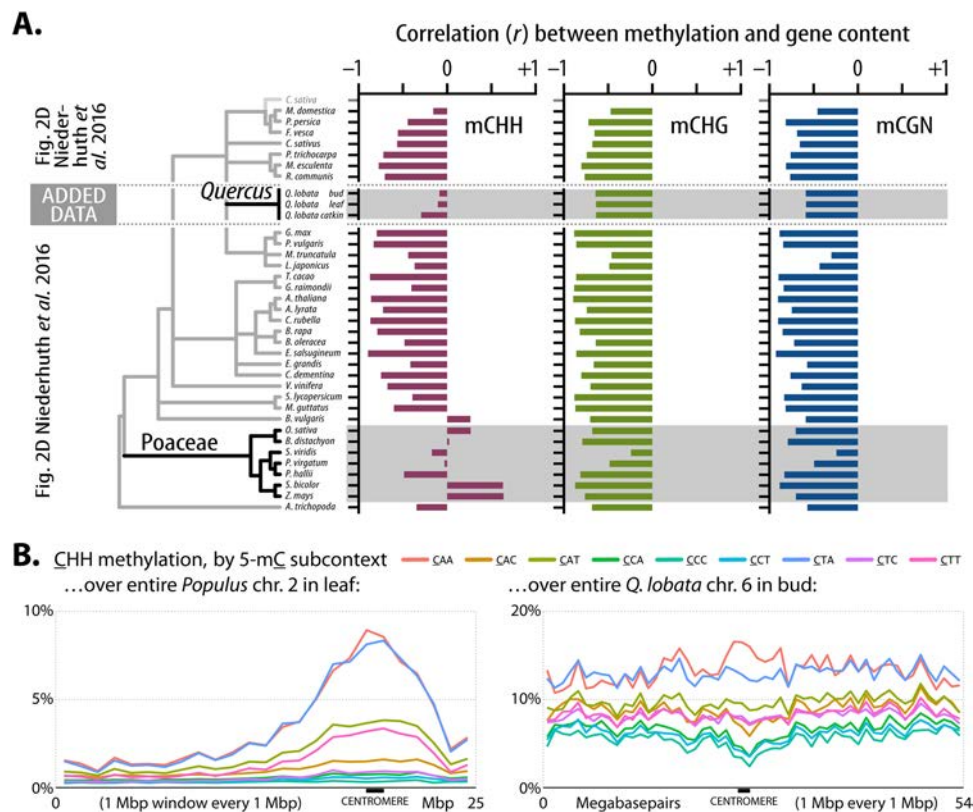
## Discussion

Our analysis of a high-contiguity, chromosome-level annotated oak genome reveals previously unreported features of oaks that

might contribute to its ability to adapt to new environments and resulting dominance in North American ecosystems. We find surprising similarities to grasses (Poaceae), another highly successful group of plants. Oaks and grasses both have genomes with large repeat-rich intergenic regions and share methylation features that are somewhat unusual, given the current limited sampling of methylomes in the literature. Interest has been growing in the adaptive potential provided by large complex intergenic regions often found in plants with larger genomes<sup>39–41</sup>. For example, a substantially higher proportion of loci associated with phenotypic variation are found in the large intergenic regions of maize versus the smaller intergenic regions of *Arabidopsis*<sup>41</sup>. Much of this regulatory variation has been found in non-TE stably unmethylated DNA<sup>42,43</sup>, such that more than 40% of phenotypic variation in maize was associated with open chromatin that makes up less than 1% of the genome<sup>44</sup>. On the other hand, the high density of diverse TEs, which has been connected with local adaptation<sup>45</sup>, can be a source of both transcription factor binding sites and regulatory non-coding RNAs<sup>46</sup>, and play a role in three-dimensional genome structure<sup>43,47,48</sup>. An abundance of intergenic heterochromatin-like structure has been demonstrated in grasses<sup>8,49,50</sup> and, based on patterns suggestive of CMT-mediated methylation<sup>37,38</sup>, are likely also found in oaks (Fig. 5I, Supplementary Figs. 19, 21 and 22, and Supplementary Data 1). Given the dramatic differences reflected in the chromosome-wide subcontext methylation patterns in the gene-dense arms of *Arabidopsis* and tomato versus the wider spread of genes in maize and rice<sup>37</sup>, and similar differences in poplar versus oak (Fig. 6B and Supplementary Fig. 23), oaks and grasses may have some regulatory strategies distinct from those in other angiosperms. Another indicator of similarity between oaks and grasses is the correlation of CHH methylation levels (mCHH) and gene count along chromosomes (Fig. 6A). A comprehensive characterization within oaks and across the angiosperms awaits further experimentation and better, more comparable genome sequences, constructed and annotated with consistent methods.

Pronounced mCHH islands are another feature shared between oaks and grasses. In maize, mCHH islands have been proposed to enforce boundaries between heterochromatin and euchromatin, and as such contribute to maintaining the suppression of TEs during increases in neighboring gene expression<sup>8,35,51</sup>. Measured as the ratio of peak mCHH to whole-genome average mCHH, we find oaks have unusually strong 5' mCHH islands (Supplementary Data 1), but it remains to be seen if they also contribute to boundary enforcement. As reported by Martin et al.<sup>52</sup> among six grass species, we found the same range of significant associations between the presence of mCHH islands and various genic features. Following their GLM model, mCHH island presence was associated with lower exonic CG methylation (estimate =  $-0.573 \pm 0.06$ ,  $z = -9.1$ ,  $p < 0.001$ ), higher gene length (estimate =  $2.95 \pm 0.33$ ,  $z = 9.1$ ,  $p < 0.001$ ), shorter distance from TSS to TE (estimate =  $-10.16 \pm 0.97$ ,  $z = -10.5$ ,  $p < 0.001$ ) and higher gene expression (estimate =  $5.30 \pm 1.76$ ,  $z = 3.0$ ,  $p < 0.01$ )<sup>52</sup>. However, our importance values were lower, and it is quite possible the mCHH islands are simply the result of the type of TEs found near gene boundaries. In valley oak (Fig. 5), maize<sup>53</sup>, and *Arabidopsis*<sup>54</sup>, mCHH is influenced by TE family, proximity to genes, and chromosomal location. The strong enrichment of small, highly methylated TE families near-genes (Fig. 5E, F) could be due to, for example, selection against large TEs in gene proximal regions.

The presence of a gene family, DUF247—known to play a role in reproductive self-recognition in rye-grass—as one of the largest and densest blocks of tandemly duplicated genes (Fig. 4A) suggests the interesting possibility of a non-self-recognition compatibility system<sup>55</sup>. Oaks have long been thought to possess some



**Fig. 6 DNA methylation across chromosomes. A** Pearson correlation ( $r$ ) between methylation level and the number of genes (100 kbp windows) for mCHH (left), mCHG (middle), and mCG (right) context levels from leaf tissue-based Fig. 2D of Niederhuth et al.<sup>36</sup>, inserting our data 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 *Populus trichocarpa*<sup>89</sup> and *Q. lobata* (and see Fig. 5 legend).

kind of self-incompatibility (SI) system because of their high outcrossing rates, but single gene SI systems have not fit observations. However, a non-self-recognition system would be consistent with observed crossing results among self, intra-, and inter-specific pollinations<sup>56</sup>. Both self- and non-self-recognition systems of co-adapted genes expressed in pollen and pistil and preventing self-fertilization have evolved independently in several lineages of angiosperms<sup>55,57</sup>. While the roles of DUF247 genes need experimental verification, their large numbers, high diversity, and conservation across species suggest some kind of collaborative system and are consistent with a non-self-recognition system that could both promote outcrossing while also permitting occasional self and interspecific crosses.

Oaks have a vast reservoir of tandemly duplicated genes of a wide variety of ages (Fig. 4F), contributing to their genetic and phenotypic diversity. As reported for pedunculate oak<sup>6</sup>, resistance genes are a particularly prominent component of the tandemly duplicated gene blocks in valley oak, especially the larger (and older) ones (see worksheets inside Supplementary Data 2 for tandem pair sizes  $\geq 20, 30, 40$ ). The three oak genomes contain hundreds to thousands of potential R-genes: *Q. lobata* has  $\approx 700$ –3000, *Q. robur* has  $\approx 600$ –2300, and *Q. suber* has  $\approx 700$ –3000. In defending oaks from bacteria, viruses, nematodes, oomycetes, and insects, R-genes may both enable the long lifespan of oaks<sup>6</sup>, and also address the puzzle of how a single or two oak species are able to dominate so many of the ecosystems they occupy. The classic Janzen–Connell ecological hypothesis proposes that pathogens promote tropical forest diversity through conspecific negative density-dependent (CNDD) mortality, but CNDD has been shown across all forest types<sup>58,59</sup>. In oaks, the high number and potential complexity of R-genes could provide a

mechanism to reduce CNDD mortality caused by pathogens<sup>60</sup>. Moreover, the large effective population size could maintain R-genes, especially if not costly<sup>61</sup>. In fact, other ecosystem-dominant trees, which also contain large numbers of domains associated with resistance genes (such as NB-ARC and LRRs), include the highly speciose *Eucalyptus* ( $\sim 600$  species) and *Populus* (Supplementary Table 7). Extensive research demonstrating the importance of R-gene diversity at both the individual and the population level is ongoing in *Arabidopsis*, crop species, and other plants<sup>62,63</sup>. Studying oaks with their large and complex pools of R-genes will provide an important extension of this work.

Inspection of our highly contiguous genome identified numerous syntenic blocks of remnant genes from the  $\gamma$  triplication event, which occurred  $\approx 120$  Mya ago when the common ancestor of angiosperms underwent two whole-genome duplication events<sup>30–32</sup>. More than 18% of protein-coding genes participate in a gene pair directly supporting a self-syntenic block (SSB), and more than a third of the genome is spanned by a manifest triplication (even without chaining blocks). SSBs (for example, Fig. 4 and Supplementary Data 2) provide an extensive single genome resource for documenting remnants associated with the  $\gamma$  event. Our annotation finds triplicate families to be enriched for transcription factors, as well as signal transduction and housekeeping genes generally (Supplementary Table 6 and Supplementary Data 2), as has been found in other studies, e.g., Rensing<sup>64</sup>. These genes, although maintained over millions of years and highly interconnected<sup>65</sup>, can respond to selective pressures by modifying their existing roles. For example, a recent study of silver birch found selective sweeps around candidate genes enriched among ancient polyploid duplicates that encode developmental timing and function in physiological cross talk<sup>7</sup>. In

oaks, it would be constructive to learn whether these ancient genes have undergone positive selection, allowing adaptation to new environments.

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

## Methods

**Study species, samples, and genomic lab work.** *Q. lobata* Née (Fagaceae) is a widely-distributed endemic California oak species found in oak savannas, oak woodlands, and riparian forests. Oak has a highly outcrossed mating system<sup>66</sup> with the potential for long-distance gene flow occurring through wind-dispersed pollen with long-tailed distributions, despite many near-neighbor pollinations<sup>67,68</sup>. Acorn dispersal is often restricted except for occasional long-distance colonization by jays<sup>69</sup>. Occupying an unglaciated region of California, contemporary populations are at least 200k years old with no evidence of severe bottlenecks during cold periods<sup>70,71</sup> like those described for European oaks from glaciation that retreated in the past 10k–20k years, allowing rapid recolonization from refugia in Italy and Spain<sup>72</sup>. Valley oak and other California oak species have been used as a reliable food source and cultural resource by native peoples of western North America for >10k years<sup>73</sup>. Since the arrival of Europeans, valley oak populations have experienced extensive habitat loss<sup>74</sup>, and current population recruitment is jeopardized by cattle grazing, rodents, and other factors<sup>75,76</sup>. Moreover, as its climate niche shifts north and upward<sup>75,77,78</sup>, extant populations are further challenged by climate warming<sup>79</sup>.

The focal tree for this study is *Q. lobata* adult SW786, which is located at the UC Santa Barbara Sedgwick Nature Reserve, and is the same individual that was sequenced for version 1.0 of the genome<sup>9</sup>. Leaf samples for the initial Illumina sequencing (532 M paired-end [PE] 250 nt reads with  $\approx 500$  nt inserts giving 133 Gnt and  $\approx 175\times$  coverage, and 318 M mate pair [MP] 150 nt reads from  $\approx 3$  knt to  $\approx 12$  knt fragments giving 48 Gnt and  $\approx 56\times$  coverage) were collected in September 2014, as described in Sork, et al.<sup>9</sup>. Additional leaves were collected and DNA extracted in April 2016 for Pacific Biosciences whole-genome SMRTbell libraries (6 M reads of mean  $\approx 9$  knt and N50  $\approx 13$  knt giving 58 Gnt and  $\approx 80\times$  coverage), and in March 2017 for Dovetail Chicago Hi-C library preparation. For details of the 19 whole-genome resequencing libraries (Illumina PE, mean  $\approx 24\times$  coverage) used for the demographic analysis, three-tissue (bud, leaf, stem) Pacific Biosciences Iso-Seq and Illumina RNA-Seq transcriptome libraries contributing to annotation, and three-tissue (bud, catkin, and young leaf) whole-genome bisulfite libraries (Illumina SE,  $\approx 18\times$ – $19\times$  coverage) for the DNA methylomes, see Supplementary Note 1. Sample collection, library preparation, sequencing, and initial data processing.

**Genome assembly.** We constructed the final genome in multiple stages. Stage 1: For the initial “Hybrid Primary” assembly (818 Mbp in 3.6k scaffolds, with longest 6.7 Mbp and NG50  $\approx 1.2$  Mbp assuming at-the-time estimated 730 Mbp for the haploid genome), we applied MaSuRCA 3.2.1<sup>80</sup> to our genomic Illumina PE, Illumina MP, and PacBio SMRT reads. The assembler identified high heterozygosity and selected diploid settings, allowing it to set aside most divergent haplotype variants; the result generally contains a single haplotype, but randomly phased, because we chose the larger scaffold whenever the assembler split two haplotypes into distinct scaffolds. Those scaffolds filtered out as alternative haplotypes were gathered into the “Hybrid Alternative” additions (466 Mbp in 17k scaffolds, with the longest 1.2 Mbp). Stage 2: To assist completeness, we aligned to Stage 1 Primary+Alternative 82k of 84k transcripts and gene fragments from a prior RNA-Seq-derived transcriptome<sup>81</sup>, with 81k aligning to Primary. To avoid loss of potential coding regions, we moved 317 scaffolds from Alternative to Primary, forming the “Hybrid-plus-Transcript Primary” assembly (872 Mbp in 4.0k scaffolds, with longest 6.7 Mbp and NG50  $\approx 1.2$  Mbp), and “Hybrid-plus-Transcript Alternative” additions (412 Mbp in 16k scaffolds, with longest 0.8 Mbp). Stage 3: We increased NG50 by aligning Stage 2 Primary scaffold ends with *bwa mem*<sup>82</sup>, merging scaffolds that had unique end matches of >94% identity longer than 40 kbp. This step created the “Hybrid-plus-Transcript-Merged Primary” assembly (861 Mbp in 3.2k scaffolds, with the longest 10.2 Mbp and NG50  $\approx 1.9$  Mbp) and “Hybrid-plus-Transcript-Merged Alternative” additions (16k scaffolds). Stage 4: Next, we generated Hi-C long-range linking information from an Illumina-sequenced library produced by Dovetail Genomics, which we used to re-scaffold with HiRise<sup>10</sup> after read alignment with a modified SNAP (<http://snap.cs.berkeley.edu>), dramatically increasing NG50. Scores from the HiRise learned likelihood model were used to identify and break presumed mis-joins,

identify prospective joins, and commit joins above a threshold; shotgun reads from Stage 1 were used to close gaps where possible. Stage 5: Finally, after HiRise, any redundant haplotype contigs remaining (that truly belong in the same place because the other haplotype in a scaffolded assembly) are expected to be adjacent to the other haplotype and this is as close as they can be placed under the linear ordering constraint of HiRise output. We used this property to remove remaining extra haplotype contigs by aligning adjacent contigs to each other and finding those smaller than their direct neighbor that had > 50% syntenic alignment with the neighbor, thereby moving 14 Mbp to Alternative and forming the final “Hi-C-Scaffolded-plus-Neighbor-Cleaned Primary” (“version 3.0”) assembly (Fig. 1C). The 12 longest scaffolds represent near full-length chromosomes (Supplementary Figs. S1 and S2) and a total of 811 Mbp (96%) of non-gap sequence.

This manuscript reports and analyzes the valley oak version 3.0 assembly. A small number of contaminants (found exclusively among non-chromosomal scaffolds) were subsequently identified and deleted, leading to the current (at time of publication) version 3.2 appearing at NCBI. These minor deletions do not significantly change our analyses or their interpretations.

**Comparisons of *Q. lobata* and *Q. robur* assemblies to linkage map.** The *Q. robur*  $\times$  *Q. petraea* linkage groups (LGs)<sup>11</sup> are taken from [http://arachne.pierroton.inra.fr/cgi-bin/cmap/map\\_set\\_info?map\\_set\\_acc=51](http://arachne.pierroton.inra.fr/cgi-bin/cmap/map_set_info?map_set_acc=51) using Supplementary Table 3 from Lepoittevin, Bodénès<sup>83</sup> as sequence-defined SNPs, dropping SNPs associated to more than one LG. Genomic locations were identified with BLASTN+ 2.2.30 ( $E < 10^{-15}$ , word size 8), keeping for each query all alignments with bitscore  $\geq 97\%$  of the top bitscore. We plot SNPs that have either a unique surviving alignment, or multiple alignments but all to the same chromosome and with chromosomal span of hits  $\leq 2$  Mbp wide.

Various alignments and visualizations appear in Figs. 1E and 4A, B, and at the project website (<https://valleyoak.ucla.edu/genomic-resources> under “Alignments and 2-D visualizations”) assisted in *Q. lobata* with discovery and identification of the tandem-like blocks of duplicated genes and self-syntenic blocks (SSBs). The principal software components employed were LASTZ for nucleotide alignments (with repeats masked by RepeatMasker after RepeatModeler); BLASTP, Diamond, and Parasail for homologous gene pair detection (on respective genome project protein-coding gene models) and subsequent detailed alignment refinement; C++, Mathematica, and Perl for scripting and pixel generation/import; and ImageMagick and Adobe Photoshop for manipulation, browsing, and annotation of generally multi-gigapixel images of high resolution (e.g., 10 kbp/pixel). Pfam hits were determined with InterProScan or direct HMMer runs.  $K_c$  was computed for all homologous protein pairs discovered with other tools by re-aligning with EMBOSS ‘needle’ (<http://emboss.sourceforge.net/>), converting to codons with ‘pal2nal.pl’ (<http://www.bork.embl.de/pal2nal/>), and finally calculating  $K_c$  with ‘codeml’ from PAML (<http://abacus.gene.ucl.ac.uk/software/paml.html>).

**Demographic history.** We inferred demographic history using the PSMC’ algorithm<sup>14</sup> by mapping the *Q. lobata* and *Q. robur* sequencing project genomic shotgun reads to their respective reference assemblies, as well as high-coverage genomic reads for 19 *Q. lobata* individuals (Supplementary Table 1) to the *Q. lobata* assembly. We called heterozygous sites in each genome (forming a VCF file with all callable sites) and composed input for PSMC’ with ‘vcfAllSiteParser.py’ (<https://github.com/stschiff/msmc-tools>). Masking and filters are as described in Supplementary Note 4. Demographic analysis—Input to PSMC’. We ran PSMC’ using default parameters except 200 for maximum number of iterations. Because PSMC’ inference can be prone to biases, we assessed robustness of our conclusions (Fig. 2B). First, we assessed if PSMC’ is capable of accurate inference by generating simulated datasets following the inferred demographic history, and re-ran inference on these (Fig. 2D). These runs suggested that the only major issue was the oldest population sizes often being over-estimated. We thus selected ancestral population sizes matching empirical genome-wide heterozygosities (Supplementary Figs. 8–10) and trimmed display in Fig. 2C accordingly (Supplementary Fig. 7 exhibits untrimmed trajectories). Next, we tested whether PSMC’ could reliably infer changes in population size on timescales relevant to *Quercus*. We simulated 10 test datasets of each run type under our presented demographic models (in Fig. 2C) using the coalescent simulator ‘msprime’<sup>84</sup>. With each simulated genome, we computed heterozygosity and used PSMC’ to infer demography; (see Supplementary Note 4. Demographic analysis—Simulations in ‘msprime’). We found accurate inference of population sizes over time, except for the single oldest time step where it tends to be over-estimated (Fig. 2D). Note that inferred demographic trajectories from whole genome-based methods such as PSMC’ can be complex but not predict empirical summary statistics such as the genome-wide distribution of heterozygosity<sup>85</sup>.

**Repetitive sequences.** The primary repeat analysis is outlined in Fig. 3A and began with the construction of a *Q. lobata*-specific database of repeat families by RepeatModeler/Classifier open-1.0.8, which was then applied with RepeatMasker open-4.0.6. Because family consensus sequences are not always full length for their class or irredundant by close sequence similarity, we applied PSI-CD-HIT 4.7 to family consensus sequences at 45% nucleotide identity (the level where, as the threshold was lowered, intracluster similarities stopped falling in frequency and

began rising) and chose a canonical rotation and strand for tandem repeat units, so as to cluster families into repeat “superfamilies” (SFs). Generally, each SF was assigned the RepeatClassifier class of the longest member of the family that was not unknown (if any; approximately two-thirds of SF-covered base pairs were classifiable). Annotated intervals for a SF are the nucleotide-level union of all intervals for member families, and SFs were assigned “s1RF####” accessions roughly serialized by descending mass. For certain uses (e.g., gene annotation), we also applied structurally-aware LTRharvest and LTRdigest from GenomeTools 1.5.9 to specifically target the abundant LTR TEs, identifying 28k instances of total mass 184 Mbp (not much larger than the 179 Mbp in LTR-classified SFs). Further details are in Supplementary Note 6. Repetitive sequences.

**Annotation of protein-coding genes.** Figure 7A outlines dataflow of the PCG modeling process employed.

**Pure Iso-Seq models.** The Pacific Biosciences of California, Inc. (PacBio) pipeline generated 197k–223k nominally full-length non-chimeric polished transcripts (reads) from the poly-A-selected strand-specific bud, leaf, and stem PacBio-sequenced Iso-Seq libraries. Pooling tissues, Minimap2 aligned each read to zero to five reference genome locations (96–99% uniquely); 11% of alignments were filtered out based on empirical criteria. Preliminary exon–intron structure was obtained by focusing on the reference side of each alignment, ignoring short insertions and merging short deletions and gapless blocks. Inspection found compact and well-isolated gene loci with generally concordant pileups at each of the 24k tentative loci, which had highly variable coverage (1–14k reads each; 56%  $\geq 5$ , 22% = 1). Inspection of higher-coverage loci found reads within pileups to vary: (i) at the exact coordinate level (with exon–intron boundaries moved by typically <10 nt vs. common); (ii) at the structural level (introns resized or deleted or inserted, generally in a small minority of reads, and at more loci and in more ways than likely by alternative splicing); and (iii) in extent (with some reads truncated, especially at the 5′ end, with loss of multiple exons possible). A consensus exon–intron model per locus was generated by resolving (i) via rounding boundaries within  $\pm 25$  nt to a most common boundary; (ii) by generally keeping exons and introns only in at least half of reads; and (iii) by extending 5′ and 3′ ends to the furthest extent observed. CDS assignments were made considering three methods: (i) longest ORFs; (ii) filtered (including restriction to only near-best hits per locus) BLASTP-equivalent (Diamond 0.9.22.123) alignments ( $E < 0.001$ ) of translations of all ORFs to the entire NCBI 2018-05-18 “nr” database (22k consensuses had at least one hit, with 83% of top hits involving at least half of both the translated ORF and the NCBI sequence, and with 99% having  $\geq 50\%$  amino acid identity, and 90% having  $E < 10^{-35}$ ; assignment of an ORF required agreement among all surviving hits); and (iii) a log<sub>2</sub>-odds bicodon coding potential trained using a selected subset of the NCBI analysis. Partial (and six-frame) ORFs were permitted. A consensus was assigned CDS (i.e., an ORF) if (ii) identified an ORF, the longest member of (i) was of the same frame with non-empty intersection with that ORF, and (iii) was also of the same frame and with non-empty intersection. This attributed CDS (and, hence, UTR5 and UTR3) to 19k loci, with 95% on the consensus read strand and 85% having  $\geq 50$  nt of both UTR5 and UTR3. Hand inspection of a random subset found them to be generally good quality (often needing no edits). Diverse AUGUSTUS hints were constructed from the Iso-Seq reads and pure Iso-Seq models for eventual use in the final AUGUSTUS run near the end of the PCG modeling process.

**AUGUSTUS bootstrap.** The 19k pure Iso-Seq models were filtered to a very high confidence subset of 2639, then thinned to 2558 by choosing single representatives from homology clusters determined via Exonerate 2.4.0 “affine:local” protein alignments. These were split uniformly at random into 1698-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 121 M to 153 M 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 was constructed from STAR observed genomic base-pair coverage and empirical splices, StringTie reference-quoted transcripts, and GMAP alignments.

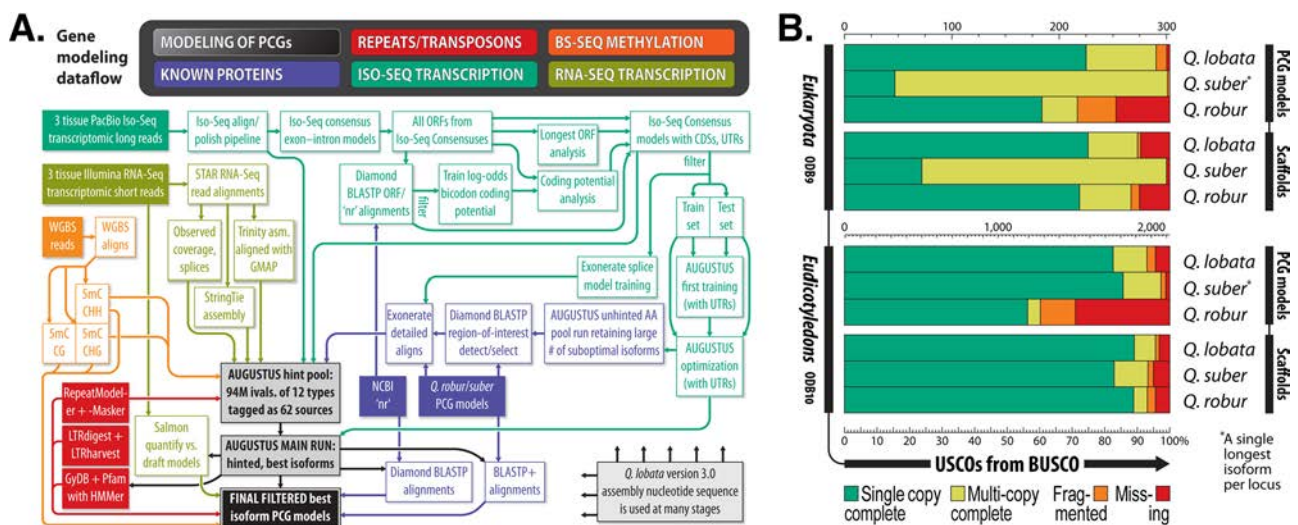
**Known proteins.** Protein translations of the *Q. robur* and *Q. suber* PCG models were BLASTP-equivalent (Diamond) aligned to a temporary trained/optimized but unhinted AUGUSTUS run generating and retaining a very large number of sub-optimal isoform models. The resultant hits were used to identify 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.

**BS-Seq DNA mCHG and mCHH 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 94 M hinting intervals of 12 types tagged as from 62 sources. (Because 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,



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

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

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

**Methylomes and analysis of tissue-specific methylation patterns.** Sample collection, library preparation, sequencing, and initial methylation calling are described in Supplementary Note 8. Methylomes and analysis of methylation patterns. Libraries were prepared using the TruSeq Nano DNA (Illumina) and Epitex kits (Qiagen), and sequenced as 100 nt single-end reads on an Illumina HiSeq 4000 to median coverage 18–19-fold. Methylation levels were determined using Methylpy v1.4.6<sup>86</sup>. DeepTools v3.1.2<sup>87</sup> “computeMatrix” and “plotProfile” were used to assess methylation levels with respect to gene models and repeat superfamilies (Fig. 5A–F, I and Supplementary Figs. 17 and 18, with default parameters except as described in legends). Methylation levels for 100 bp windows were calculated by dividing the total number of reads calling “T” (= methylated) by the total number of informative reads (“C” or “T”) for all genomic cytosine positions in the appropriate sequence context within the window. Genome-wide average methylation levels (Fig. 5 and Supplementary Figs. 16 and 18) were calculated by averaging 100 bp window levels for the 12 chromosomal scaffolds. Per-site methylation levels in Supplementary Fig. 16 were calculated by dividing reads calling methylation (“T”) by all informative reads (“C” or “T”) for each position, plotting with R “ggplot2” v3.3.2<sup>88</sup>. Designation of genomic regions with respect to genes (1 kbp up, 5’ UTR, etc.) was done with Bedtools v2.27.1 (<https://doi.org/10.1093/bioinformatics/btq033>) and “bed12toAnnotation.awk” (<https://github.com/guigolab/geneid/blob/master/scripts/bed12toAnnotation.awk>). PCG model spans do not overlap in our annotation; however, overlaps for 1 kbp upstream and 1 kbp downstream regions were removed from the 1k up and 1k down categories, including overlaps that spanned neighboring genes. Gene regions overlapping with intervals (200 kbp to 3 Mbp) covering pericentromeric regions were removed. Introns were separated into first intron vs. other introns. Chromosome scale plots of subcontext methylation (Figs. 5I and 6A) were calculated via Bedtools as the mean of the percent methylation at each genomic cytosine position in the appropriate sequence context within each 1 Mb window, every 1 Mb. *Populus* methylation data<sup>89</sup> was for Tree 13 branch 1, obtained from GEO (<https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE132939> 2020). Local correlations between methylation levels and gene count were determined using methods as close as possible to Niederhuth, Bewick<sup>36</sup> to maximize the relevance of the comparison. Thus, using Bedtools, the genome was divided into 100 kbp windows with 50 kbp overlaps. Methylation for each 100 kbp window was from averaging 100 bp window methylation levels (as above). Genes per window were counted with Bedtools “intersect,” requiring at least 50% of the gene span to be inside the window. Pearson correlation between gene count and methylation level was calculated with R “cor” with incomplete observations dropped.

**Reporting summary.** Further information on research design is available in the Nature Research Reporting Summary linked to this article.

## Data availability

The version 3.0 assembly and associated data are available variously at NCBI under Umbrella BioProject PRJNA781973, the project website (<https://valleyoak.ucla.edu/>), and the project UCSC genome browser (<http://genomes.mcdb.ucla.edu/cgi-bin/hgTracks?db=queLob3>). DNA sequencing reads for this project are available from SRA under accession numbers SRX1088964–SRX10889720 (75 files) and SRX10972769–SRX10972854 (86 files).

## Code availability

The code for PSMC files and masks is available at <https://doi.org/10.5281/zenodo.5899535><sup>90</sup>.

Received: 12 April 2021; Accepted: 11 March 2022;

Published online: 19 April 2022

## References

- Kremer, A. & Hipp, A. L. Oaks: an evolutionary success story. *N. Phytol.* **226**, 987–1011 (2020).
- Hipp, A. L., Manos, P. S. & Cavender-Bares, J. How oak trees evolved to rule the forests of the Northern Hemisphere. *Sci. Am.* **323**, 42–49 (2020).
- Barrón, E. et al. in *Oaks Physiological Ecology. Exploring the Functional Diversity of Genus Quercus L.* (eds Gil-Pelegrín, E. P.-P. J. & Sancho-Knapik, D.) 39–105 (Springer, 2017).
- Denk, T., Grimm, G. W., Manos, P. S., Deng, M. & Hipp, A. L. in *Oaks Physiological Ecology. Exploring the Functional Diversity of Genus Quercus L.* (eds Gil-Pelegrín, E., Peguero-Pina, J. & Sancho-Knapik, D.) 13–38 (Springer, 2017).
- Cavender-Bares, J. Diversity, distribution, and ecosystem services of the North American oaks. *J. Int. Oak Soc.* **27**, 37–48 (2016).
- Plomion, C. et al. Oak genome reveals facets of long lifespan. *Nat. Plants* **4**, 440–452 (2018).
- Salojärvi, J. et al. Genome sequencing and population genomic analyses provide insights into the adaptive landscape of silver birch. *Nat. Genet.* **49**, 904–912 (2017).
- Li, Q. et al. RNA-directed DNA methylation enforces boundaries between heterochromatin and euchromatin in the maize genome. *Proc. Natl Acad. Sci. USA* **112**, 14728–14733 (2015).
- Sork, V. L. et al. First draft assembly and annotation of the genome of a California endemic oak. *Quercus lobata* Née (Fagaceae). *G3 Genes Genomes Genet.* **11**, 3485–3495 (2016).
- Putnam, N. H. et al. Chromosome-scale shotgun assembly using an in vitro method for long-range linkage. *Genome Res.* **26**, 342–350 (2016).
- Bodénès, C., Chancerel, E., Ehrenmann, F., Kremer, A. & Plomion, C. High-density linkage mapping and distribution of segregation distortion regions in the oak genome. *DNA Res.* **23**, 115–124 (2016).
- Ramos, A. M. et al. The draft genome sequence of cork oak. *Sci. Data* **5**, 180069 (2018).
- Hipp, A. L. et al. Genomic landscape of the global oak phylogeny. *N. Phytol.* **226**, 1198–1212 (2020).
- Schiffels, S. & Durbin, R. Inferring human population size and separation history from multiple genome sequences. *Nat. Genet.* **46**, 919–925 (2014).
- Rundel, P. W. et al. Mediterranean biomes: evolution of their vegetation, floras, and climate. *Annu. Rev. Ecol. Syst.* **47**, 383–407 (2016).
- Corbett-Detig, R. B., Hartl, D. L. & Sackton, T. B. Natural selection constrains neutral diversity across a wide range of species. *PLoS Biol.* **13**, e1002112 (2015).
- Myburg, A. A. et al. The genome of *Eucalyptus grandis*. *Nature* **510**, 356–362 (2014).
- Argout, X. et al. The genome of *Theobroma cacao*. *Nat. Genet.* **43**, 101–108 (2011).
- Denoued, F. et al. The coffee genome provides insight into the convergent evolution of caffeine biosynthesis. *Science* **345**, 1181–1184 (2014).
- Tuskan, G. A. et al. The genome of black cottonwood, *Populus trichocarpa* (Torr. & Gray). *Science* **313**, 1596–1604 (2006).
- Stanke, M., Schöffmann, O., Morgenstern, B. & Waack, S. Gene prediction in eukaryotes with a generalized hidden Markov model that uses hints from external sources. *BMC Bioinformatics* **7**, 62 (2006).
- El-Gebali, S. et al. The Pfam protein families database in 2019. *Nucleic Acids Res.* **47**, D427–D432 (2018).
- Seppely, M., Manni, M. & Zdobnov, E. M. in *Gene Prediction. Methods in Molecular Biology* (ed. Kollmar, M.) 227–245 (Humana, 2019).
- Gururani, M. A. et al. Plant disease resistance genes: current status and future directions. *Physiol. Mol. Plant Pathol.* **78**, 51–65 (2012).
- Mistry, J. et al. Pfam: The protein families database in 2021. *Nucleic Acids Res.* **49**, D412–D419 (2020).
- Manzanares, C. et al. A gene encoding a DUF247 domain protein cosegregates with the S Self-Incompatibility locus in perennial ryegrass. *Mol. Biol. Evol.* **33**, 870–884 (2015).
- Harkess, A. et al. The asparagus genome sheds light on the origin and evolution of a young Y chromosome. *Nat. Commun.* **8**, 1279 (2017).
- Harkness, A. & Brandvain, Y. Non-self recognition-based self-incompatibility can alternatively promote or prevent introgression. *N. Phytol.* **231**, 1630–1643 (2021).
- Veckman, E. et al. Overcoming challenges in variant calling: exploring sequence diversity in candidate genes for plant development in perennial ryegrass (*Lolium perenne*). *DNA Res.* **26**, 1–12 (2019).
- Chanderbali, A. S., Berger, B. A., Howarth, D. G., Soltis, D. E. & Soltis, P. S. Evolution of floral diversity: genomics, genes and gamma. *Philos. Trans. R. Soc. B Biol. Sci.* <https://doi.org/10.1098/rstb.2015.0509> (2017).
- Jiao, Y. N. et al. A genome triplication associated with early diversification of the core eudicots. *Genome Biol.* **13**, R3 (2012).
- Vekemans, D. et al. Gamma paleohexaploidy in the stem lineage of core eudicots: significance for MADS-Box gene and species diversification. *Mol. Biol. Evol.* **29**, 3793–3806 (2012).
- Birchler, J. A. & Veitia, R. A. The gene balance hypothesis: implications for gene regulation, quantitative traits and evolution. *N. Phytol.* **186**, 54–62 (2010).
- Higo, A. et al. DNA methylation is reconfigured at the onset of reproduction in rice shoot apical meristem. *Nat. Commun.* **11**, 4079–4079 (2020).
- Gent, J. I. et al. CHH islands: de novo DNA methylation in near-gene chromatin regulation in maize. *Genome Res.* **23**, 628–637 (2013).
- Niederhuth, C. E. et al. Widespread natural variation of DNA methylation within angiosperms. *Genome Biol.* **17**, 194 (2016).

37. Gouil, Q. & Baulcombe, D. C. DNA methylation signatures of the plant chromomethyltransferases. *PLoS Genet.* **12**, e1006526 (2016).
38. Song, X. & Cao, X. Context and complexity: analyzing methylation in trinucleotide sequences. *Trends Plant Sci.* **22**, 351–353 (2017).
39. Carpentier, M.-C. et al. Retrotranspositional landscape of Asian rice revealed by 3000 genomes. *Nat. Commun.* **10**, 24 (2019).
40. Choi, J. Y. & Lee, Y. C. G. Double-edged sword: the evolutionary consequences of the epigenetic silencing of transposable elements. *PLoS Genet.* **16**, e1008872 (2020).
41. Mei, W., Stetter, M. G., Gates, D. J., Stitzer, M. C. & Ross-Ibarra, J. Adaptation in plant genomes: bigger is different. *Am. J. Bot.* **105**, 16–19 (2018).
42. Crisp, P. A. et al. Stable unmethylated DNA demarcates expressed genes and their cis-regulatory space in plant genomes. *Proc. Natl Acad. Sci. USA* **117**, 23991–24000 (2020).
43. Ricci, W. A. et al. Widespread long-range cis-regulatory elements in the maize genome. *Nat. Plants* **5**, 1237–1249 (2019).
44. Rodgers-Melnick, E., Vera, D. L., Bass, H. W. & Buckler, E. S. Open chromatin reveals the functional maize genome. *Proc. Natl Acad. Sci. USA* **113**, E3177–E3184 (2016).
45. Baduel, P., Quadrana, L., Hunter, B., Bomblies, K. & Colot, V. Relaxed purifying selection in autopolyploids drives transposable element over-accumulation which provides variants for local adaptation. *Nat. Commun.* **10**, 5818 (2019).
46. Chuong, E. B., Elde, N. C. & Feschotte, C. Regulatory activities of transposable elements: from conflicts to benefits. *Nat. Rev. Genet.* **18**, 71–86 (2017).
47. Li, E. et al. Long-range interactions between proximal and distal regulatory regions in maize. *Nat. Commun.* **10**, 2633 (2019).
48. Vanrobays, E., Thomas, M. & Tatout, C. Heterochromatin positioning and nuclear architecture. *Ann. Plant Rev. online* pp. 157–190 (2018).
49. Makarevitch, I. et al. Genomic distribution of maize facultative heterochromatin marked by trimethylation of H3K27. *Plant Cell* **25**, 780–793 (2013).
50. Zhao, L. et al. Chromatin loops associated with active genes and heterochromatin shape rice genome architecture for transcriptional regulation. *Nat. Commun.* **10**, 3640 (2019).
51. Long, J. C. et al. Decrease in DNA methylation 1 (DDM1) is required for the formation of mCHH islands in maize. *J. Integr. Plant Biol.* **61**, 749–764 (2019).
52. Martin, G. T., Seymour, D. K. & Gaut, B. S. CHH methylation islands: a nonconserved feature of grass genomes that is positively associated with transposable elements but negatively associated with gene-body methylation. *Genome Biol. Evol.* **13**, evab144 (2021).
53. Achour, Z. et al. Low temperature triggers genome-wide hypermethylation of transposable elements and centromeres in maize. Preprint at *bioRxiv* <https://doi.org/10.1101/573915> (2019).
54. Sasaki, E., Kawakatsu, T., Ecker, J. R. & Nordborg, M. Common alleles of CMT2 and NRPE1 are major determinants of CHH methylation variation in *Arabidopsis thaliana*. *PLoS Genet.* **15**, e1008492 (2020).
55. Iwano, M. & Takayama, S. Self/non-self discrimination in angiosperm self-incompatibility. *Curr. Opin. Plant Biol.* **15**, 78–83 (2012).
56. Boavida, L. C., Silva, J. P. & Feijo, J. A. Sexual reproduction in the cork oak (*Quercus sober* L.) - II. Crossing intra- and interspecific barriers. *Sex. Plant Reprod.* **14**, 143–152 (2001).
57. Charlesworth, D., Vekemans, X., Castric, V. & Glémin, S. Plant self-incompatibility systems: a molecular evolutionary perspective. *N. Phytol.* **168**, 61–69 (2005).
58. Johnson, D. J., Beaulieu, W. T., Bever, J. D. & Clay, K. Conspecific negative density dependence and forest diversity. *Science* **336**, 904–907 (2012).
59. Bever, J. D., Mangan, S. A. & Alexander, H. M. Maintenance of plant species diversity by pathogens. *Annu. Rev. Ecol. Evol. Syst.* **46**, 305–325 (2015).
60. Marden, J. H. et al. Ecological genomics of tropical trees: how local population size and allelic diversity of resistance genes relate to immune responses, cosusceptibility to pathogens, and negative density dependence. *Mol. Ecol.* **26**, 2498–2513 (2017).
61. Stump, S. M., Marden, J. H., Beckman, N. G., Mangan, S. A. & Comita, L. S. Resistance genes affect how pathogens maintain plant abundance and diversity. *Am. Nat.* **196**, 472–486 (2020).
62. Xue, J.-Y., Takken, F. L. W., Nepal, M. P., Maekawa, T. & Shao, Z.-Q. Editorial: Evolution and functional mechanisms of plant disease resistance. *Front. Genet.* <https://doi.org/10.3389/fgene.2020.593240> (2020).
63. Karasov, T. L., Shirsekar, G., Schwab, R. & Weigel, D. What natural variation can teach us about resistance durability. *Curr. Opin. Plant Biol.* **56**, 89–98 (2020).
64. Rensing, S. A. Gene duplication as a driver of plant morphogenetic evolution. *Curr. Opin. Plant Biol.* **17**, 43–48 (2014).
65. Defoort, J., Van de Peer, Y. & Carretero-Paulet, L. The evolution of gene duplicates in angiosperms and the impact of protein–protein interactions and the mechanism of duplication. *Genome Biol. Evol.* **11**, 2292–2305 (2019).
66. Sork, V., Dyer, R., Davis, F. & Smouse, P. Mating system in California Valley oak, *Quercus lobata* Neé. In: *Proc. Fifth Symposium on Oak Woodlands: Oaks in California's Changing Landscape* (eds Standiford, R. B. M. D. & Purcell, K. L.) 427–444 (Pacific Southwest Research Station, Forest Service, U.S. Department of Agriculture, 2002).
67. Pluess, A. R. et al. Short distance pollen movement in a wind-pollinated tree, *Quercus lobata* (Fagaceae). *For. Ecol. Manag.* **258**, 735–744 (2009).
68. Sork, V. L. & Smouse, P. E. Genetic analysis of landscape connectivity in tree populations. *Landsc. Ecol.* **21**, 821–836 (2006).
69. Sork, V. L., Smouse, P. E., Grivet, D. & Scofield, D. G. Impact of asymmetric male and female gamete dispersal on allelic diversity and spatial genetic structure in valley oak (*Quercus lobata* Née). *Evol. Ecol.* **29**, 927–945 (2015).
70. Grivet, D., Deguilloux, M.-F., Petit, R. J. & Sork, V. L. Contrasting patterns of historical colonization in white oaks (*Quercus* spp.) in California and Europe. *Mol. Ecol.* **15**, 4085–4093 (2006).
71. Gugger, P. F., Ikegami, M. & Sork, V. L. Influence of late Quaternary climate change on present patterns of genetic variation in valley oak, *Quercus lobata* Née. *Mol. Ecol.* **22**, 3598–3612 (2013).
72. Petit, R. J. et al. Chloroplast DNA footprints of postglacial recolonization by oaks. *Proc. Natl Acad. Sci. USA* **94**, 9996–10001 (1997).
73. Anderson, M. K. *Tending the Wild: Native American Knowledge and the Management of California's Natural Resources* (University of California Press, 2005).
74. Whipple, A. A., Grossinger, R. M. & Davis, F. W. Shifting baselines in a California oak savanna: Nineteenth century data to inform restoration scenarios. *Restor. Ecol.* **19**, 88–101 (2011).
75. McLaughlin, B. C. & Zavaleta, E. S. Predicting species responses to climate change: demography and climate microrefugia in California valley oak (*Quercus lobata*). *Glob. Change Biol.* **18**, 2301–2312 (2012).
76. Tyler, C. M., Kuhn, B. & Davis, F. W. Demography and recruitment limitations of three oak species in California. *Q. Rev. Biol.* **81**, 127–152 (2006).
77. Sork, V. L. et al. Gene movement and genetic association with regional climate gradients in California valley oak (*Quercus lobata* Née) in the face of climate change. *Mol. Ecol.* **19**, 3806–3823 (2010).
78. Kueppers, L. N., Snyder, M. A., Sloan, L. C., Zavaleta, E. S. & Fulfrost, B. Modeled regional climate change and California endemic oak ranges. *Proc. Natl Acad. Sci. USA* **102**, 16281–16286 (2005).
79. Browne, L., Wright, J. W., Fitz-Gibbon, S., Gugger, P. F. & Sork, V. L. Adaptational lag to temperature in valley oak (*Quercus lobata*) can be mitigated by genome-informed assisted gene flow. *Proc. Natl Acad. Sci. USA* **50**, 25179–25185 (2019).
80. Zimin, A. V. et al. Hybrid assembly of the large and highly repetitive genome of *Aegilops tauschii*, a progenitor of bread wheat, with the MaSuRCA megareads algorithm. *Genome Res.* **27**, 787–792 (2017).
81. Cokus, S. J., Gugger, P. F. & Sork, V. L. Evolutionary insights from *de novo* transcriptome assembly and SNP discovery in California white oaks. *BMC Genomics* **16**, 552 (2015).
82. Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics* **25**, 1754–1760 (2009).
83. Lepoittevin, C. et al. Single-nucleotide polymorphism discovery and validation in high-density SNP array for genetic analysis in European white oaks. *Mol. Ecol. Resour.* **15**, 1446–1459 (2015).
84. Kelleher, J., Etheridge, A. M. & McVean, G. Efficient coalescent simulation and genealogical analysis for large sample sizes. *PLoS Comput. Biol.* **12**, e1004842 (2016).
85. Beichman, A. C., Phung, T. N. & Lohmueller, K. E. Comparison of single genome and allele frequency data reveals discordant demographic histories. *G3 Genes Genomes Genet.* **7**, 3605–3620 (2017).
86. Schultz, M. D. et al. Human body epigenome maps reveal noncanonical DNA methylation variation. *Nature* **523**, 212–216 (2015).
87. Ramírez, F. et al. deepTools2: a next generation web server for deep-sequencing data analysis. *Nucleic Acids Res.* **44**, W160–W165 (2016).
88. Wickham, H. *ggplot2: Elegant Graphics for Data Analysis* (Springer-Velgag, 2016).
89. Hofmeister, B. T. et al. A genome assembly and the somatic genetic and epigenetic mutation rate in a wild long-lived perennial *Populus trichocarpa*. *Genome Biol.* **21**, 259 (2020).
90. Garcia, J., Zhen, Y. & Lohmueller, K. Demographic history analysis scripts for *Quercus lobata* reference genome (v1.0.2). Zenodo <https://doi.org/10.5281/zenodo.5899420> (2022).
91. Griffin, J. R. & Critchfield, W. B. *The Distribution of the Forest Trees in California* (Pacific SW Forest and Range Experiment Station, U.S. Department of Agriculture Forest Service, 1972).

## Acknowledgements

We acknowledge the native peoples of California whose relationship with the land has allowed the peoples of California to benefit from oak ecosystems, and allowed our research team to study oaks, and especially the Chumash as the traditional land caretakers of the land that is home of tree SW786. We acknowledge the UC-Santa Barbara Sedgwick Reserve and the University of California Natural Reserve System. We thank Krista Beckley, Dylan Burge,

and Andy Lentz for fieldwork; Krista Beckley for DNA extractions and library preparation; Marco Morselli for assistance with bisulfite sequencing library preparation; Suhua Feng for assistance with DNA and RNA sequencing; and Luke Browne for maps in Fig. 2. Sequencing was conducted at the following core facilities: Illumina and PacBio were carried out at the DNA Technologies and Expression Analysis Cores of the UC Davis Genome Center, supported by NIH Shared Instrumentation Grant 1S10OD010786-01; and whole-genome resequencing, WGBS, and RNA-Seq for demographics, methylation, and transcriptomic studies, respectively, were done at the UCLA Broad Stem Cell Genome Core facility. We thank Lily Shiue and Thomas Swale, Dovetail Genomics, for the high-quality scaffolding performed by HiRise. This research was supported by the National Science Foundation Plant Genome Research Program (NSF Award IOS-1444661).

### Author contributions

V.L.S., M.P., and S.L.S. conceived the overall project design and management and obtained grant support; V.L.S. initiated, coordinated, and supervised the project and manuscript. S.J.C. annotated and analyzed genes and repeats, designed and conducted genome comparative analyses, and created/wrote figures, results, methods, and Supplementary Information (SI) for these sections. S.T.F.-G. analyzed methylomes, designed and conducted comparative methylation analyses, created/wrote figures, results, discussion, methods, and SI for this topic; called genetic variants (GATK) for demographics; and submitted genomic resources for public availability. A.V.Z. and D.P. assembled and curated the genome sequence and contributed to results, methods, and SI for these sections. J.G. and Y.Z. analyzed genetic variation data. J.G. conducted the demographic analysis. K.E.L. designed and supervised the demographic analysis and J.G. and K.E.L. contributed results, methods, and SI for this section. S.T.F.-G. and S.J.C. examined DUF247. C.L.H. conducted lab preparation for DNA sequencing, resequencing, bisulfite sequencing, and RNA sequencing. V.L.S., S.J.C., S.T.F.-G., A.V.Z., J.G., P.F.G., K.E.L., M.P., and S.L.S. edited text. S.J.C. edited manuscript figures. S.J.C., S.T.F.-G., and V.L.S. interpreted data and wrote the manuscript.

### Competing interests

The authors declare no competing interests.

### Additional information

**Supplementary information** The online version contains supplementary material available at <https://doi.org/10.1038/s41467-022-29584-y>.

**Correspondence** and requests for materials should be addressed to Victoria L. Sork.

**Peer review information** *Nature Communications* thanks John Mackay, Chad Niederhuth, and the other anonymous reviewer(s) for their contribution to the peer review of this work.

**Reprints and permission information** is available at <http://www.nature.com/reprints>

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2022