

## Reporting Summary

Nature Research wishes to improve the reproducibility of the work that we publish. This form provides structure for consistency and transparency in reporting. For further information on Nature Research policies, see [Authors & Referees](#) and the [Editorial Policy Checklist](#).

### Statistical parameters

When statistical analyses are reported, confirm that the following items are present in the relevant location (e.g. figure legend, table legend, main text, or Methods section).

n/a Confirmed

- The exact sample size ( $n$ ) for each experimental group/condition, given as a discrete number and unit of measurement
- An indication of whether measurements were taken from distinct samples or whether the same sample was measured repeatedly
- The statistical test(s) used AND whether they are one- or two-sided  
*Only common tests should be described solely by name; describe more complex techniques in the Methods section.*
- A description of all covariates tested
- A description of any assumptions or corrections, such as tests of normality and adjustment for multiple comparisons
- A full description of the statistics including central tendency (e.g. means) or other basic estimates (e.g. regression coefficient) AND variation (e.g. standard deviation) or associated estimates of uncertainty (e.g. confidence intervals)
- For null hypothesis testing, the test statistic (e.g.  $F$ ,  $t$ ,  $r$ ) with confidence intervals, effect sizes, degrees of freedom and  $P$  value noted  
*Give  $P$  values as exact values whenever suitable.*
- For Bayesian analysis, information on the choice of priors and Markov chain Monte Carlo settings
- For hierarchical and complex designs, identification of the appropriate level for tests and full reporting of outcomes
- Estimates of effect sizes (e.g. Cohen's  $d$ , Pearson's  $r$ ), indicating how they were calculated
- Clearly defined error bars  
*State explicitly what error bars represent (e.g. SD, SE, CI)*

Our web collection on [statistics for biologists](#) may be useful.

### Software and code

Policy information about [availability of computer code](#)

Data collection

1/ sequencing data collection: we used Roche/454 (non operating anymore) and Illumina sequencing machines that included their own computer packages for DNA sequence collection. 2/ genotyping data collection: Allele calling from the MassArray iPLEX assay was processed in Typer Viewer v 4.0.26.75 software (Agena Bioscience).

Data analysis

custom-made scripts were made available through public web sites. This is clearly stated in the Ms.  
We relied on R for statistical analyses.  
The following public or commercial software were used: AMAS / Augustus / BLAST / Bowtie2 / bcftools / BWA-MEM / BLAT / bedtools / BioNJ / Celera assembler / COANCESTRY / CAFE version 3.1 / Censor / CLC genomics Workbench / Clustal-Omega / Cabog / ClustalW / DUST / EuGene / EggLib / FGENESH / FigTree V1.4.3 / Fasttree 2.1.8 / FEELnc version 26.05.2015 / FeatureCounts / Ggplot2 / G-block v0.91b / Geneious 6.1.8 / HAPLOMERGE V1 / HMMER 3.0 / Interproscan v5.13-52.0 / infernal 1.1 / Jellyfish / LPmerge / LTRharvest / MAP / MuTect / MUSCLE / MUMmer / MAFFT version 5 / MEGA5 / MEGA6 / Newbler assembler (version MapAsmResearch-04/19/2010-patch-08/17/2010) / netGen2 / NUCmer / NBLR parser / OLC assembler / OrthoMCL / Picard / Popoolation2 / Popoolation / PAML 4 / PhyML 3.0 / Prank / RepeatMasker open-3.0/ RepeatScout / REPET / RAXML 7.7.2 / RNAmmer / SSPACE / Sickle / SNPA gene finder / SpliceMachine / Signl P / SAMtools / STRUCTURE / SOAPdenovo2 / SiLiX / S-MART / Seaview version 4 / Stringtie v1.0.1 / TRF / TargetP / TMHMM / Typer Viewer v4.0.26.75 / topGO 2.22.0 / tandem repeatsFinder / trimAl gt 0.2 / STAR 2.4.0i / tRNAscan-SE / taxize / TreeDyn v198.3 / YASS

For manuscripts utilizing custom algorithms or software that are central to the research but not yet described in published literature, software must be made available to editors/reviewers upon request. We strongly encourage code deposition in a community repository (e.g. GitHub). See the Nature Research [guidelines for submitting code & software](#) for further information.

## Data

Policy information about [availability of data](#)

All manuscripts must include a [data availability statement](#). This statement should provide the following information, where applicable:

- Accession codes, unique identifiers, or web links for publicly available datasets
- A list of figures that have associated raw data
- A description of any restrictions on data availability

The oak haploid genome assembly and corresponding annotation have been deposited in the European Nucleotide Archive under project accession code PRJEB19898. Other sequence release data are indicated in Supplementary tables 1, 13, 14 and 19 and Supplementary Data Set 10. We also invite readers to download data stored at the URLs indicated in section 6 (Web resources) as well as in the oakgenome web site: <http://www.oakgenome.fr>.

## Field-specific reporting

Please select the best fit for your research. If you are not sure, read the appropriate sections before making your selection.

Life sciences       Behavioural & social sciences

For a reference copy of the document with all sections, see [nature.com/authors/policies/ReportingSummary-flat.pdf](https://nature.com/authors/policies/ReportingSummary-flat.pdf)

## Life sciences

### Study design

All studies must disclose on these points even when the disclosure is negative.

Sample size	selection of tree and herbaceous genomes for comparison of gene content were made according to the quality of publicly available genom assemblies: see Online Methods section "Detection of significant expansion/contraction in woody perennials. "
Data exclusions	As explained in the Method section (line 711-713) we excluded from our initial population screening : related genotypes (3) as well as introgressed genotypes (8), leaving as much as possible unrelated and not introgressed genotypes for population genetics analysis.
Replication	It is very important to stress that somatic mutations we deemed reliable were detected by comparing multiple sequencing libraries, taking into account the chronology of branch development. This is a much more powerful validation than would have been provided by the technical Sanger validation. In our approach, polymorphisms generated by assembly errors would be very unlikely i/ to show differences among the libraries for the different branches, and ii/ to follow a temporal pattern coherent with the chronology of branch development. Regarding Sanger validation, according to a recent study in humans (Beck et al. 2016, <a href="https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4878677/">https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4878677/</a> ), Sanger sequencing was more likely to incorrectly refute a true positive variant from NGS than to correctly identify a false positive. In addition, Sanger sequencing appears to have low sensitivity for low frequency variants, as expected for the vast majority of somatic mutations. We therefore preferred to rely on our comparative approach.
Randomization	We compared estimates for genetic between genes from expanded, contracted, and unchanged gene families (orthogroups) in oak. We accounted for the different gene family sizes, by randomly sampling 1000 genes from each of these three categories and repeating the operation 100 times.
Blinding	genotype calls obtained from the Typer Viewersoftware of the mass array spectrometer were controlled by two people (according to a complete double-blind design), and considered as valid when the two calls were identical.

## Materials & experimental systems

Policy information about [availability of materials](#)

n/a	Involvement in the study
<input type="checkbox"/>	<input checked="" type="checkbox"/> Unique materials
<input checked="" type="checkbox"/>	<input type="checkbox"/> Antibodies
<input checked="" type="checkbox"/>	<input type="checkbox"/> Eukaryotic cell lines
<input checked="" type="checkbox"/>	<input type="checkbox"/> Research animals
<input checked="" type="checkbox"/>	<input type="checkbox"/> Human research participants

### Unique materials

Obtaining unique materials	the genotype selected to sequence the oak genome is a 100 year-old tree in our experimntal station. Still living and plant material are available upon request.
----------------------------	---

## Antibodies

Antibodies used

Describe all antibodies used in the study; as applicable, provide supplier name, catalog number, clone name, and lot number.

Validation

Describe the validation of each primary antibody for the species and application, noting any validation statements on the manufacturer's website, relevant citations, antibody profiles in online databases, or data provided in the manuscript.

## Eukaryotic cell lines

Policy information about [cell lines](#)

Cell line source(s)

State the source of each cell line used.

Authentication

Describe the authentication procedures for each cell line used OR declare that none of the cell lines used were authenticated.

Mycoplasma contamination

Confirm that all cell lines tested negative for mycoplasma contamination OR describe the results of the testing for mycoplasma contamination OR declare that the cell lines were not tested for mycoplasma contamination.

Commonly misidentified lines  
(See [ICLAC](#) register)

Name any commonly misidentified cell lines used in the study and provide a rationale for their use.

## Research animals

Policy information about [studies involving animals](#); [ARRIVE guidelines](#) recommended for reporting animal research

Animals/animal-derived materials

For laboratory animals, report species, strain, sex and age OR for animals observed in or captured from the field, report species, sex and age where possible.

## Human research participants

Policy information about [studies involving human research participants](#)

Population characteristics

Describe the covariate-relevant population characteristics of the human research participants (e.g. age, gender, genotypic information, past and current diagnosis and treatment categories).

# Method-specific reporting

n/a	Included in the study
<input checked="" type="checkbox"/>	<input type="checkbox"/> ChIP-seq
<input checked="" type="checkbox"/>	<input type="checkbox"/> Flow cytometry
<input checked="" type="checkbox"/>	<input type="checkbox"/> Magnetic resonance imaging

## ChIP-seq

Data deposition

Confirm that both raw and final processed data have been deposited in a public database such as [GEO](#).

Confirm that you have deposited or provided access to graph files (e.g. BED files) for the called peaks.

Data access links

May remain private before publication.

For "Initial submission" or "Revised version" documents, provide reviewer access links. For your "Final submission" document, provide a link to the deposited data.

Files in database submission

Provide a list of all files available in the database submission.

Genome browser session

(e.g. [UCSC](#))

Provide a link to an anonymized genome browser session for "Initial submission" and "Revised version" documents only, to enable peer review. Write "no longer applicable" for "Final submission" documents.

## Methodology

Replicates

Describe the experimental replicates, specifying number, type and replicate agreement.

Sequencing depth

Describe the sequencing depth for each experiment, providing the total number of reads, uniquely mapped reads, length of reads and whether they were paired- or single-end.

Antibodies

Describe the antibodies used for the ChIP-seq experiments; as applicable, provide supplier name, catalog number, clone name, and lot number.

Peak calling parameters

Specify the command line program and parameters used for read mapping and peak calling, including the ChIP, control and index files used.

Data quality

Describe the methods used to ensure data quality in full detail, including how many peaks are at FDR 5% and above 5-fold

Data quality	<i>enrichment.</i>
Software	<i>Describe the software used to collect and analyze the ChIP-seq data. For custom code that has been deposited into a community repository, provide accession details.</i>

## Flow Cytometry

### Plots

Confirm that:

- The axis labels state the marker and fluorochrome used (e.g. CD4-FITC).
- The axis scales are clearly visible. Include numbers along axes only for bottom left plot of group (a 'group' is an analysis of identical markers).
- All plots are contour plots with outliers or pseudocolor plots.
- A numerical value for number of cells or percentage (with statistics) is provided.

### Methodology

Sample preparation	<i>Describe the sample preparation, detailing the biological source of the cells and any tissue processing steps used.</i>
Instrument	<i>Identify the instrument used for data collection, specifying make and model number.</i>
Software	<i>Describe the software used to collect and analyze the flow cytometry data. For custom code that has been deposited into a community repository, provide accession details.</i>
Cell population abundance	<i>Describe the abundance of the relevant cell populations within post-sort fractions, providing details on the purity of the samples and how it was determined.</i>
Gating strategy	<i>Describe the gating strategy used for all relevant experiments, specifying the preliminary FSC/SSC gates of the starting cell population, indicating where boundaries between "positive" and "negative" staining cell populations are defined.</i>

Tick this box to confirm that a figure exemplifying the gating strategy is provided in the Supplementary Information.

## Magnetic resonance imaging

### Experimental design

Design type	<i>Indicate task or resting state; event-related or block design.</i>
Design specifications	<i>Specify the number of blocks, trials or experimental units per session and/or subject, and specify the length of each trial or block (if trials are blocked) and interval between trials.</i>
Behavioral performance measures	<i>State number and/or type of variables recorded (e.g. correct button press, response time) and what statistics were used to establish that the subjects were performing the task as expected (e.g. mean, range, and/or standard deviation across subjects).</i>

### Acquisition

Imaging type(s)	<i>Specify: functional, structural, diffusion, perfusion.</i>
Field strength	<i>Specify in Tesla</i>
Sequence & imaging parameters	<i>Specify the pulse sequence type (gradient echo, spin echo, etc.), imaging type (EPI, spiral, etc.), field of view, matrix size, slice thickness, orientation and TE/TR/flip angle.</i>
Area of acquisition	<i>State whether a whole brain scan was used OR define the area of acquisition, describing how the region was determined.</i>
Diffusion MRI	<input type="checkbox"/> Used <input type="checkbox"/> Not used

### Preprocessing

Preprocessing software	<i>Provide detail on software version and revision number and on specific parameters (model/functions, brain extraction, segmentation, smoothing kernel size, etc.).</i>
Normalization	<i>If data were normalized/standardized, describe the approach(es): specify linear or non-linear and define image types used for transformation OR indicate that data were not normalized and explain rationale for lack of normalization.</i>
Normalization template	<i>Describe the template used for normalization/transformation, specifying subject space or group standardized space (e.g. original Talairach, MNI305, ICBM152) OR indicate that the data were not normalized.</i>

## Noise and artifact removal

*Describe your procedure(s) for artifact and structured noise removal, specifying motion parameters, tissue signals and physiological signals (heart rate, respiration).*

## Volume censoring

*Define your software and/or method and criteria for volume censoring, and state the extent of such censoring.*

## Statistical modeling &amp; inference

## Model type and settings

*Specify type (mass univariate, multivariate, RSA, predictive, etc.) and describe essential details of the model at the first and second levels (e.g. fixed, random or mixed effects; drift or auto-correlation).*

## Effect(s) tested

*Define precise effect in terms of the task or stimulus conditions instead of psychological concepts and indicate whether ANOVA or factorial designs were used.*

Specify type of analysis:  Whole brain  ROI-based  Both

Statistic type for inference  
(See [Eklund et al. 2016](#))

*Specify voxel-wise or cluster-wise and report all relevant parameters for cluster-wise methods.*

## Correction

*Describe the type of correction and how it is obtained for multiple comparisons (e.g. FWE, FDR, permutation or Monte Carlo).*

## Models &amp; analysis

n/a | Involved in the study

- Functional and/or effective connectivity  
  Graph analysis  
  Multivariate modeling or predictive analysis

## Functional and/or effective connectivity

*Report the measures of dependence used and the model details (e.g. Pearson correlation, partial correlation, mutual information).*

## Graph analysis

*Report the dependent variable and connectivity measure, specifying weighted graph or binarized graph, subject- or group-level, and the global and/or node summaries used (e.g. clustering coefficient, efficiency, etc.).*

## Multivariate modeling and predictive analysis

*Specify independent variables, features extraction and dimension reduction, model, training and evaluation metrics.*

## Behavioural &amp; social sciences

## Study design

All studies must disclose on these points even when the disclosure is negative.

## Study description

*Briefly describe the study type including whether data are quantitative, qualitative, or mixed-methods (e.g. qualitative cross-sectional, quantitative experimental, mixed-methods case study).*

## Research sample

*State the research sample (e.g. Harvard university undergraduates, villagers in rural India) and provide relevant demographic information (e.g. age, sex) and indicate whether the sample is representative. Provide a rationale for the study sample chosen. For studies involving existing datasets, please describe the dataset and source.*

## Sampling strategy

*Describe the sampling procedure (e.g. random, snowball, stratified, convenience). Describe the statistical methods that were used to predetermine sample size OR if no sample-size calculation was performed, describe how sample sizes were chosen and provide a rationale for why these sample sizes are sufficient. For qualitative data, please indicate whether data saturation was considered, and what criteria were used to decide that no further sampling was needed.*

## Data collection

*Provide details about the data collection procedure, including the instruments or devices used to record the data (e.g. pen and paper, computer, eye tracker, video or audio equipment) whether anyone was present besides the participant(s) and the researcher, and whether the researcher was blind to experimental condition and/or the study hypothesis during data collection.*

## Timing

*Indicate the start and stop dates of data collection. If there is a gap between collection periods, state the dates for each sample cohort.*

## Data exclusions

*If no data were excluded from the analyses, state so OR if data were excluded, provide the exact number of exclusions and the rationale behind them, indicating whether exclusion criteria were pre-established.*

## Non-participation

*State how many participants dropped out/declined participation and the reason(s) given OR provide response rate OR state that no participants dropped out/declined participation.*

## Randomization

*If participants were not allocated into experimental groups, state so OR describe how participants were allocated to groups, and if*

