1     **Comparative Genomics of Six *Juglans* Species Reveals Disease-associated Gene**
2                          **Family Contractions.**
3
4

5     Alexander Trouern-Trend[1]*, Taylor Falk[1]*, Sumaira Zaman[1], Madison Caballero[1], David
6     B. Neale[2], Charles H. Langley[4], Abhaya Dandekar[2], Kristian A. Stevens[3,4+], Jill L.
7     Wegrzyn[1+]
8
9     [1]Department of Ecology and Evolutionary Biology, University of Connecticut, Storrs,
10    CT, USA
11    [2]Department of Plant Sciences, University of California Davis, Davis, CA, USA
12    [3]Department of Computer Science, University of California Davis, Davis, CA, USA
13    [4]Department of Evolution and Ecology, University of California Davis, Davis, CA, USA
14
15    *Joint First Authors
16    +Joint Corresponding Authors (jill.wegrzyn@uconn.edu, kastevens@ucdavis.edu)
17

23

1

24 **ABSTRACT**
25 *Juglans* (walnuts), the most speciose genus in the walnut family (Juglandaceae)
26 represents most of the family's commercially valuable fruit and wood-producing trees.
27 It includes several species used as rootstock in agriculture for their resistance to various
28 abiotic and biotic stressors. We present the full structural and functional genome
29 annotations of six *Juglans* species and one outgroup within Juglandaceae (*Juglans regia, J.*
30 *cathayensis, J. hindsii, J. microcarpa, J. nigra, J. sigillata* and *Pterocarya stenoptera*) produced
31 using BRAKER2 semi-unsupervised gene prediction pipeline and additional tools. For
32 each annotation, gene predictors were trained using 19 tissue-specific *J. regia*
33 transcriptomes aligned to the genomes. Additional functional evidence and filters were
34 applied to multi-exonic and mono-exonic putative genes to yield between 27,000 and
35 44,000 high-confidence gene models per species. Comparison of gene models to the
36 BUSCO embryophyta dataset suggested that, on average, genome annotation
37 completeness was 85.6%. We utilized these high-quality annotations to assess gene
38 family evolution within *Juglans* and among *Juglans* and selected Eurosid species. We
39 found notable contractions in several gene families in *J. hindsii*, including disease
40 resistance-related Wall-associated Kinase (WAK) and *Catharanthus roseus* Receptor-like
41 Kinase (CrRLK1L) and others involved in abiotic stress response. Finally, we confirmed
42 an ancient whole genome duplication that took place in a common ancestor of
43 Juglandaceae using site substitution comparative analysis.
44
45 **INTRODUCTION**
46 It is anticipated that new genomic resources for *Juglans* (walnuts) will lead to improved
47 timber and nut production by accelerating development of advanced agriculture,
48 breeding efforts, and resource management techniques for the genus. Already, these
49 practices are beneficiaries of the genomic analyses and tool development potentiated by
50 the growing pool of *Juglans* sequence data (Martínez-García *et al.* 2017; Bernard *et al.*
51 2018; Marrano *et al.* 2018; Famula *et al.* 2019; Zhu *et al.* 2019). The recent publication of
52 the unannotated draft reference genomes of six *Juglans* species: *J. nigra* (Eastern black
53 walnut), *J. hindsii* (Hinds black walnut), *J. microcarpa* (Texas walnut), *J. sigillata* (iron
54 walnut), *J. cathayensis* (Chinese walnut), *J. regia* (Persian or English walnut) and a
55 member of the sister group to *Juglans, Pterocarya stenoptera* (Chinese wingnut) greatly
56 expands the existing resource and provides an unprecedented opportunity to apply
57 tools such as genomic selection to the Juglandaceae (Stevens *et al.* 2018).
58
59 The genomes considered are Eurasian and North American species from across the
60 three sections of the genus, *Cardiocaryon, Dioscaryon* (syn. sect. *Juglans*) and *Rhysocaryon*
61 (Figure 1). These species were selected for their historical and agricultural importance
62 and for their phylogenetic placements, which span the breadth of the genus. The North
63 American species *J. nigra, J. hindsii,* and *J. microcarpa* are members of sect. *Rhysocaryon*

1

64      and grow in riparian forests in the eastern, western and southern United States,
65      respectively. *J. hindsii* and *J. microcarpa* are distributed in moderately dispersed
66      populations. Conversely, *J. nigra's* range is more contiguous and expansive and
67      overlaps with the northeastern distribution of *J. microcarpa*. Among these, *J. nigra* is
68      especially valued for its high-quality timber, cold-hardiness, and disease resistance
69      (Beineke 1983; Settle *et al.* 2015; Chakraborty *et al.* 2015; Chakraborty *et al.* 2016). *J. hindsii*
70      is characterized as vigorous, drought tolerant, and resistant to *Armillaria* root rot (honey
71      fungus) (Buzo *et al.* 2009).  The need for disease-resistant and climate-tolerant cultivars
72      for nut production has driven the use of *J. nigra*, *J. hindsii* and *J. microcarpa* in
73      hybridization trials (Browne *et al.* 2015).  These species contribute resistance to many
74      soilborne pathogens, including *Phytophthora* and a range of nematodes.
75
76      The sampled Eurasian species include *J. regia*, the predominant nut producer of all
77      cultivated *Juglans*. *J. regia's* western native range continues into Southeastern Europe, an
78      impressive relic of silk road trading that underlines its agricultural suitability
79      (Pollegioni *et al.* 2015). *J. sigillata* is phylogenetically adjacent to *J. regia* as the only other
80      species of sect. *Dioscaryon*, and, like *J. regia*, is valued for both its timber and nut
81      (Weckerle *et al.* 2005). Despite their divergent morphology, which includes number of
82      leaflets and nut characteristics, a growing body of molecular evidence suggests that *J.
83      regia* and *J. sigillata* may not qualify as separate species (Gunn *et al.* 2010; Zhao *et al.*
84      2018). *J. sigillata* grows sympatrically with *J. regia* and *J. cathayensis* (sect. *Cardiocaryon*) in
85      southwestern China, but the distribution of *J. cathayensis* extends beyond the sympatric
86      zone several hundred kilometers eastward to the coastline and northward towards the
87      Gobi desert. *J. cathayensis* is endangered in its natural range in China (Zhang *et al.* 2015)
88      and is evaluated in breeding programs for its resistance to lesion nematodes.  The
89      outgroup, *P. stenoptera* is also native to southeastern China and is used for ornamental
90      planting, timber, and medicinal extracts.  *P. stenoptera* has been integrated into
91      hybridization trials as a non-viable (inconsistent grafting) rootstock for its resistance to
92      *Phytophthora* (Browne *et al.* 2011).
93
94      Stevens *et. al* (2018) annotated a set of microsyntenic regions containing polyphenol
95      oxidase loci to confirm a gene duplication in an ancestral *Juglans*. Here, we describe the
96      full gene annotations of these diploid genomes, a critical missing component to their
97      full utilization as a genomic resource. We demonstrate their utility by investigating the
98      genomes in a comparative manner. First, we view the evolution of gene families in
99      *Juglans* across the phylogeny. Second, we leverage the annotations for a comparative
100     genomic analysis to date an ancient whole genome duplication in an ancestral species of
101     Juglandaceae.
102
103     **RESULTS**

104      *Semi-unsupervised Gene Prediction*
105      Errors introduced from genome annotation often lead to inconsistent gene expression
106      estimates and contribute to the inaccurate characterization of gene space, gene family
107      evolution and timing of whole genome duplications (Vijay *et al.* 2013, Denton *et al.*
108      2014). Our approach was applied across all seven genomes that leveraged RNA-Seq
109      reads generated from tissue-specific libraries of *J. regia* (Table 1). This approach took
110      advantage of the deep sequencing by directly aligning reads to the genome to resolve
111      challenges associated with reliance on error-prone and often fragmented *de novo*
112      assembled transcripts (Hoff *et al.* 2016). The bias introduced by using RNA-Seq reads
113      solely from *J. regia* for the annotation of all genomes was partially mitigated by the
114      semi-supervised training of the gene prediction tool, AUGUSTUS, included in
115      BRAKER. The AUGUSTUS component utilizes the evidence of successful alignments to
116      learn features of the genome in question and propose gene models. Repeat libraries
117      were generated and subsequently used for masking between roughly 44% and 48% of
118      the genomes prior to read alignment (Table 2; Table S1). The raw reads aligned across
119      the genomes at rates inversely proportional to their phylogenetic distance from *J. regia*.
120      Average alignment rates across *J. regia* transcriptome libraries are displayed as total
121      mapped (concordantly mapped): 87.1% (84.2%) in *J. regia*, 88.3% (78.7%) in *J. sigillata*,
122      51.8% (49.1%) in *J. cathayensis*, 68.0% (49.0%) in *J. nigra*, 41.0% (38.7%) in *J. microcarpa*,
123      49.7% (48.0%) in *J. hindsii* and 33.2% (31.2%) in the outgroup, *P. stenoptera*. Initial gene
124      prediction estimates from BRAKER2 ranged from 81,753 (*J. hindsii*) to 133,963 (*P.*
125      *stenoptera*) (Table S2). Filtering of BRAKER2 models considered completeness (start and
126      stop codon present), isoforms, exon lengths, intron lengths, and splice sites. These
127      considerations provided a reduced set of gene models for each genome: 82,610 in *J.*
128      *nigra*, 76,847 in *J. hindsii*, 114,573 in *J. microcarpa*, 83,457 in *J. sigillata*, 97,312 in *J.*
129      *cathayensis*, 84,098 in *J. regia*, and 123,420 in *P. stenoptera* (Table S2).
130
131      *Functional Annotation Filtering of Gene Models*
132      Functional annotation via sequence similarity search (SSS) and gene family assignment
133      (GFA) provides a form of validation for the proposed models and assesses their
134      completeness. The reciprocal style search required query and target sequence coverage
135      to pass a set threshold, which helped to eliminate unlikely models and validate the final
136      models. A total of 29,046 models with both SSS result and GFA, and 5,190 with only
137      GFA composed the final *J. nigra* set. The same approach was used for the 29,382 models
138      in *J. hindsii*, 43,051 in *J. microcarpa*, 27,596 in *J. sigillata*, 34,857 in *J. cathayensis*, 31,621 in *J.*
139      *regia*, and 45,808 in *P. stenoptera* (Table 2; Table S2). Structural assessment of the genes
140      examined splice sites, exons per gene, CDS lengths, and intron lengths (Table 2) and
141      reported average gene/CDS lengths relative to other angiosperm species. The vast
142      majority (> 98% in all species) of the splice sites were canonical (GT/AG). All other
143      splice site detected were GC/AG variants.

144
145 *Benchmarking Genome Annotation Completeness*
146 The embryophyta collection of 1440 single copy orthologs derived from OrthoDB can be
147 accessed via BUSCO to estimate the completeness of a plant genome assembly,
148 transcriptome, or set of gene models. These 1440 genes (Embryophyta *odb9*) were
149 aligned to all Juglandaceae assemblies and final gene models (Figure 2; Table S3, Table
150 S4). Across members of *Juglans* genus, BUSCO identified 87 to 93% of their database
151 when evaluated against the genome (Table S5). When provided with filtered complete
152 (full-length) gene models, BUSCO reported 77 to 86% completeness, and 78 to 89% with
153 partial (5′ or 3′ complete) models (Table S4).
154
155 *Orthologous Group Construction*
156 Two OrthoFinder analysis that differed in adjacent clade inclusivity were used to
157 identify homology relationships between genes of the selected species. One run of
158 annotated Juglandaceae (6 *Juglans*, 1 *Pterocarya*) species, and another including
159 previously annotated genomes from across the Eurosid superorder (13 species, see
160 Methods). OrthoFinder assigned 216,778 (92.5%) of the 234,455 total genes from the
161 Juglandaceae set to 26,458 orthogroups (Figure 3B, File S1). The resulting orthogroups
162 range in size from 2 to 190 genes. A total of 161 genes (0.1%) are in 56 species-specific
163 orthogroups. Of the *Juglans* species, *J. cathayensis* had the most genes designated to
164 species-specific orthogroups (24 genes in 8 orthogroups). Just over half, 14,429
165 orthogroups, have gene membership from all species. A total of 661 orthogroups (5268
166 genes) are represented by all *Juglans* species (excluding *Pterocarya*). The Juglandaceae
167 set included 538 orthogroups (1159 genes) specific to *Juglans* sect. *Dioscaryon* and 437
168 orthogroups (1608 genes) specific to members of *Juglans* sect. *Rhysocaryon*. Within
169 *Rhysocaryon*, 905 genes formed 389 orthogroups specific to the parapatric species, *J.*
170 *microcarpa* and *J. nigra*, but not found in the geographically isolated *J. hindsii*. A total of
171 149 orthogroups (564 genes) were specific to the three Eurasian *Juglans* species (*J. regia*,
172 *J. sigillata* and *J. cathayensis*). Genes that could not be assigned to orthogroups, included:
173 2025 (6.6%) in *J. regia*, 2801 (8.2%) in *J. cathayensis*, 1138 (4.0%) in *J. hindsii*, 3181 (7.6%) in
174 *J. microcarpa*, 934 (3.3%) in *J. nigra*, 1251 (4.7%) in *J. sigillata*, and 6347 (14.3%) in *P.*
175 *stenoptera* (Figure 3A; Figure 3B; File S1).
176
177 OrthoFinder analysis of selected Eurosid species assigned 401,186 (92.8.%) of the
178 456,424 total genes to 22,189 orthogroups (File S2). Of these, a total of 3054 genes (0.7%)
179 are present in 488 species-specific orthogroups and 6722 orthogroups contained at least
180 one gene from each species. The addition of peripheral species to the analysis resulted
181 in an increased gene contribution per species in the orthogroups. This trend is reflected
182 by fewer orthogroups resulting from the Eurosid clustering and the approximate

4

183     halving of the number of unassigned *Juglans* genes in the Eurosid clustering when
184     compared to the Juglandaceae clustering (Table S6).
185
186     *Analysis of Gene Family Evolution*
187     An evaluation of gene families among the annotated species was successful in detecting
188     significant changes between taxa. Prior to gene family analysis with CAFE, orthogroups
189     were filtered to exclude large families (> 100 gene copies) and those composed entirely
190     of paralogs.  This removed 57 of 26,458 (0.2%) orthogroups from the Juglandaceae set,
191     and 1880 of 22,189 (8.5%) orthogroups from the Eurosid set. Calculated lambda values
192     were 0.02396 and 0.02197 for Juglandaceae and Eurosid sets, respectively. The higher
193     lambda of Juglandaceae set indicates a higher calculated average rate of gene family
194     evolution. Of the 460 significant rapidly evolving orthogroups discovered based on the
195     Eurosid set, 153 (+131 families expanded/-22 families contracted) had significant
196     changes in *J. microcarpa*, 102 (+57/-45) in *J. regia*, 86 (+62/-24) in *J. cathayensis*, 76 (+30/-46)
197     in *J. sigillata*, 61 (+22/-39) in *J. nigra*, 58 (+32/-26) in *J. hindsii*, and 139 (+113/-24) in *P.*
198     *stenoptera* (Figure 5, File S4). The Juglandaceae set revealed 430 significant rapidly
199     evolving gene families of which 168 (+123/-45) had significant size changes in the *J.*
200     *microcarpa* terminal branch, 141 (+86/-55) in *J. regia*, 92 (+72/-20) in *J. cathayensis*, 98 (+39/-
201     59) in *J. sigillata*, 101 (+32/-69) in *J. nigra*, 77 (+40/-37) in *J. hindsii*, and 98 (+67/-31) in *P.*
202     *stenoptera* (File S3).
203
204     *Rhysocaryon Gene Family Evolution*
205     At the ancestral *Rhysocaryon* node, 4 significant expansions and 2 significant
206     contractions were discovered. Functional annotation of Juglandaceae orthogroups
207     expanded in *J. microcarpa* revealed high incidence of transferase activity (GO:0016740)
208     which occurred in 8 of 123 orthogroup annotations. An orthogroup annotated as
209     ankyrin repeat-containing (OG0000093) was significantly expanded in both *J. microcarpa*
210     (22 genes) and *J. sigillata* (16 genes) relative to other species (0-6 genes). Three
211     orthogroups annotated as Kinesin-like protein KIN-4C, Phosphatidylinositol 4-kinase
212     gamma 7 (P4KG7) and RNA-dependent RNA polymerase (OG0002363, OG0001584 and
213     OG0013906) were expanded in *J. microcarpa* and *J. nigra* relative to other annotated
214     species. An activating signal cointegrator orthogroup (OG0022144) with a zinc finger-
215     C2HC5 (Pfam:PF06221) domain was expanded (+13) in *J. microcarpa*. OG0000386,
216     annotated as topless-related protein 1 (TPR1) was also expanded (+6). Three
217     orthogroups annotated as "wall-associated receptor kinase-like" (OG0000502,
218     OG0000046 and OG0000685) lacked gene models from both *J. hindsii* and *J. microcarpa*.
219     OG0000685 also lacked *J. sigillata* gene models. A Heat Shock Cognate 70 kDa (HSC70)
220     orthogroup (OG0000060) was expanded in both *J. hindsii* (23 genes) and *J. microcarpa* (30
221     genes) relative to all species outside of *Rhysocaryon* (1-2 genes) and unexpectedly lacked
222     gene models from *J. nigra*. Similarly, SAPK10-like serine/threonine kinase orthogroup

5

223  (OG0001146) was also expanded in both *J. hindsii* (8 genes) and *J. microcarpa* (7 genes)
224  relative to other species (0-4 genes) and lacked *J. nigra* gene models.
225
226  The large Juglandaceae callose synthase 3-like orthogroup (OG0000004) is absent in *J.*
227  *nigra* and highly contracted in *J. microcarpa* and *J. cathayensis* (7 genes) relative to other
228  species (32-36 genes). Four cyclic nucleotide-gated ion channel orthogroups involved in
229  Plant-pathogen interaction (KEGG:04626), are lost or highly contracted in *J. nigra*:
230  OG0000038 (-7), OG0000567 (-3), OG0000177 (-4), OG0000603 (-2). Juglandaceae
231  REDUCED WALL ACETYLATION 2 (RWA2) (OG0000145), putative disease resistance
232  protein (OG0000022) and receptor-like protein kinase FERONIA-like (OG0000471)
233  orthogroups lacked *J. hindsii* gene models despite being represented by every other
234  species.
235
236  *Dioscaryon Gene Family Evolution*
237  At the ancestral *Dioscaryon* node, 5 significant expansions and 8 significant contractions
238  were discovered. Annotated gene family expansions specific to *Dioscaryon* include (+3)
239  ABC transporter B family orthogroup (OG0000286), and (+2) Ethanolamine-phosphate
240  (OG0001347) orthogroups. Juglandaceae cationic peroxidase (CEVI16) orthogroup
241  (OG0000173) related to Phenylpropanoid biosynthesis (GO:0009699) is contracted in
242  *Juglans* sect. *Dioscaryon*. Probable reticuline oxidase families (OG0000562, OG0000531)
243  annotated as containing BBE (Pfam:PF08031) and FAD binding 4 (Pfam:PF1565) lack
244  *Dioscaryon* gene models while all non-*Dioscaryon* species contribute at least 3 gene
245  copies in each orthogroup. *Dioscaryon* gene models were absent in nodulin-like
246  orthogroup (OG0000206) (-3 genes). Contractions in F-box protein orthoroup
247  (OG0000054) and Oxygen-evolving enhancer protein 2 (OG0000122) were also observed
248  (-4 and -3 genes, respectively). One SWIM zinc finger orthogroup (OG0000510) lacked
249  gene models in *J. regia, J. sigillata* and *J. microcarpa*. Another orthogroup (OG0000266)
250  annotated as SWIM zinc finger appeared to also be absent in *J. regia, J. sigillata* and *J.*
251  *microcarpa*, but a *J. sigillata* ortholog was discovered as a loss through the absence of
252  protein to genome alignment.
253
254  Gene family expansions in *J. regia* include (+4) 26s proteasome regulatory subunit
255  (OG0019963), (+3) thaumatin-like protein (OG0000263), (+4) STOMATAL
256  CYTOKINESIS DEFECTIVE 1-like (OG0001205), (+3) mitogen-activated protein kinase
257  kinase kinase (OG0012715), (+4) Hydroxyproline O-galactosyltransferase GALT6
258  (OG0004422), (+10) tubulin beta-6 chain (OG0000238). Expanded orthogroups in *J.*
259  *sigillata* include (+11) endoribonuclease dicer (OG0000131).
260
261  *Gene Family Evolution Enrichment*
262  EggNOG gene descriptions of rapidly evolving gene families were examined to infer

263      the major functional categories of rapidly expanding and contracting gene families
264      across Juglandaceae. Of the 333 instances of gene family contraction calculated across
265      the Juglandaceae, the most frequent GO molecular function terms, included: 26
266      transferase activity (GO:0016740), 9 lyase activity (GO:0016829), and 9 cyclase activity
267      (GO:0009975) families. High occurrence EggNOG-derived gene family descriptions of
268      contracting orthogroups included 25 that contained "resistance", 54 containing
269      "kinase", 8 "cytochrome P450" and 7 "channel". For the 428 instances of gene family
270      expansion, the most frequent molecular function annotations were 29 transferase
271      activity (GO:0016740), 8 transmembrane transporter activity (GO:0022857) and 7
272      heterocyclic compound binding (GO:1901363). High occurrence EggNOG descriptions
273      of expanding orthogroups include 51 containing "kinase", 19 that contained
274      "resistance" 17 that contained "synthase". Comparisons of annotated rapidly evolving
275      gene families among Juglandaceae species did reveal disproportionate gains and losses.
276      *J. microcarpa*, for example has 7 instances of expansion in "synthase" orthogroups while
277      *J. sigillata* has 0 and *J. hindsii* demonstrates contraction of 10 "kinase" orthogroups,
278      while only 2 such contractions were calculated in *J. cathayensis* (0 at the preceding node
279      shared with *Dioscaryon*). These divergent patterns of gene family evolution underline
280      the importance of having comprehensive genetic resources for multiple species within a
281      single clade. The six *Juglans* genome annotations provide an immediate reference for
282      one another and construct a genetic background for the genus.
283
284      Of the 153 significant gene family size changes in *J. microcarpa*, 131 represent
285      expansions. The changes in other *Juglans* species are more evenly distributed between
286      expansions and contractions. The inflated number of significant expansions in *J.*
287      *microcarpa* likely reflects uncollapsed heterozygosity left behind by the genome
288      assembly process, especially given the unexpectedly large size of the *J. microcarpa*
289      assembly (Table 1). A similar, but less pronounced pattern is observed in *J. cathayensis*.
290      *Selection Analysis*
291      The likelihoods of one-ratio (null), nearly neutral (NN) and positive selection (PS)
292      models were compared (Table S7). Of the 15 gene families that were tested, the nearly
293      neutral model fit the data significantly better than the null for 2 orthogroups and the
294      positive selection model for 10. Of these, 6 orthogroups (OG0000038, OG0000567,
295      OG0000603, OG0001146, OG0001205, OG0001222) were found to be under positive
296      selection across the selected sequences. OG0000038 (PS $\omega = 1.67$), OG0000567 (PS $\omega =$
297      1.76) and OG0000603 (PS $\omega = 1.84$) were annotated as cyclic nucleotide-gated ion channel
298      proteins, OG0001146 (PS $\omega = 1.24$) as a serine threonine-protein kinase, OG0001205 (PS
299      $\omega = 9.96$) annotated as STOMATAL CYTOKINESIS DEFECTIVE 1-like and OG0001222
300      (PS $\omega = 2.42$) as Chitinase-3.
301
302      *Divergence Estimates*

303    We estimated the distribution of nucleotide substitution rates at silent codon positions
304    between each of the *Juglans* genomes studied and the outgroup *Pterocarya stenoptera*. For
305    each pairwise analysis, we observed similar bimodal distributions of synonymous
306    substitution rates (Ks) between syntenic blocks of genes (Figure 4A). For these syntenic
307    blocks of genes, a whole genome duplication event would give rise to such a bimodal
308    distribution in time to the most recent common ancestor. For each species pair, we thus
309    estimated the two modes of the distribution (Table S8). The estimates for the higher
310    mode ranged from a low of Ks = 0.356 to a high of Ks = 0.364 with an average value of
311    Ks = 0.361. The lower mode ranged from a low of Ks = 0.050 to a high of Ks = 0.054 with
312    an average value of Ks = 0.053. While the non-synonymous substitution rates (Kn)
313    between syntenic blocks of genes were much lower, the distributions were also bi-
314    modal in appearance (Figure 4B)
315
316    The annotation of the genome of *Quercus robur* (oakgenome.fr) allowed us to perform
317    the same analysis with a species whose common ancestor predates the whole genome
318    duplication event common to the Juglandaceae. We chose the genomes of *J. regia* and *P.*
319    *stenoptera* as the best representatives of their genera. In both cases, while the histogram
320    was much sparser due to the additional divergence, a single prominent peak was
321    observed. For *J. regia* against *Q. robur*, it was observed at a value of Ks = 0.49 and for *P.*
322    *stenoptera* against *Q. robur*, it was observed at a value of Ks = 0.53. These divergence
323    estimates are greater than all values estimated in the *Juglans-Pterocarya* comparisons.
324
325    **DISCUSSION**
326    In this study, we utilized a comprehensive *J. regia* transcriptome dataset to produce
327    high-quality genome annotations of six recently assembled species within *Juglans* and a
328    single member of the sister genus, *Pterocarya*. The gene model set completeness as
329    measured by BUSCO suggests our annotation pipeline is suitable for comprehensive
330    capture of protein-coding genes. It is still expected that limitations of single species
331    RNA-Seq as the training input introduced some bias in the annotations for the other
332    Juglandaceae. Although the gene prediction software, BRAKER2 seems to return far
333    fewer false positive gene models than alternative applications, the process of removing
334    the extraneous models remains essential to producing genome annotations that can be
335    leveraged by the community. Still, complex plant genomes, especially those derived
336    from short read dominant assemblies, remain challenging to annotate and existing
337    pipelines typically introduce errors and false positives (Van Bel et al. 2019). The gene
338    model filtration steps presented here handled multi-exonic and mono-exonic genes
339    separately and examined both structural and functional qualities of models to permit
340    only those of the highest confidence. This phylogenetically comprehensive set of
341    diploid genome annotations represents an invaluable resource for comparative
342    genomics studies within *Juglans* and for other clades (Tuskan *et al.* 2018).

8

343
344    The species annotated in this study represent each of the three sections of *Juglans*
345    (*Cardiocaryon, Juglans* (syn. *Dioscaryon*) and *Rhysocaryon*) and represent fully the
346    diversity in the genus. These annotations will serve as a platform for identifying genetic
347    underpinnings of high-value agricultural characteristics such as drought tolerance and
348    disease resistance that are scattered across the various species (Bernard *et al.* 2018).
349    Moreover, they have the potential to add a new dimension to the ongoing medicinal
350    natural products search within *Juglans* (Yao *et al.* 2012; Xu *et al.* 2013; Kim *et al.* 2018).
351
352    Because this dataset is representative of the diversity in *Juglans*, it allows for exceptional
353    resolution of patterns in gene family evolution. Multiple samples within sect.
354    *Rhysocaryon* and sect. *Dioscaryon* increase confidence that observed patterns across those
355    genomes are true and not artifacts of technological and biological challenges.
356
357    *Challenges in Assessing Gene Family Evolution*
358    Given the nature of short read assemblies, the possibility of an assembly or annotation
359    error resulting in an incorrect consensus and falsely 'pseudogenizing' a gene model is
360    non-zero. These errors, especially in small gene families, could be interpreted as
361    significant contractions in the CAFE analysis. The weighty consequence of this effect on
362    interpreting gene family evolution underscores the importance of deep sequencing for
363    comparative studies, and as the shift towards long read sequencing progresses,
364    adherence to best base-calling and polishing practices.
365
366    The risk of introducing false positive expansions is most prominent in the genome
367    assembly phase. High heterozygosity in parts of a genome make the recovery of both
368    haplotypes (for diploids) difficult for those regions. In final assemblies the haplotypes
369    are often reported in separate contigs. Any gene models prevailing in these regions will
370    falsely occur in duplicate within the annotation if the haplotigs are not recognized. The
371    *J. microcarpa, P. stenoptera*, and to a lesser extent, *J. cathayensis* genomes exhibited these
372    patterns by showing high duplication rates in BUSCO analyses (Table S4), inflated
373    numbers of gene models (Table S2), and larger than expected genome sizes (Table 1).
374    The evidence for uncollapsed heterozygosity in these genomes was reinforced by the
375    absence of an additional peak representing taxa-specific duplications in the Ks
376    distributions. Computational tools have been developed to address the challenges of
377    resolving heterozygous region but are most effective when applied to long-read (or
378    hybrid) assemblies (Chin et al. 2016).
379
380    The vastly reduced cost of sequencing over the past several years has enabled genus-
381    level analysis of whole genome diversity, a scale at which it becomes tractable to assess
382    patterns and significance of changes in gene CNV and other structural variation. Given

383 the newness of this capability, a sharp increase in sequencing projects capable of
384 resolving CNV should be expected. However, there are still only a few studies that have
385 established the phenotypic and fitness consequences of CNV (Cook *et al.* 2012,
386 Würschum *et al.* 2018) and even fewer that involve full-genome assessments (Prunier *et*
387 *al.* 2018).
388
389 Convergent shifts in copy number under strong selective pressure for glyphosate
390 resistance were reported for the *EPSPS* gene in eight weedy species (Patterson *et al.*
391 2018). This finding is notable because it points towards modulated gene expression
392 levels through CNV as a potential source of rapid adaptation on short timescales. These
393 types of structural variations most often occur in genomic regions called CNV hotspots,
394 which are enriched for low-copy repeats (LCRs) (Hastings *et al.* 2009). In a genome-wide
395 survey, distinguishing between an ancestral event and parallel evolution would require
396 attention to the entire duplicated genomic region in each taxon. These investigations
397 lend a greater importance to the production of near chromosome-level assemblies
398 because poor contiguity obscures the ability to resolve structural variants.
399
400 A recent pangenome study in *Poplar* showed that intraspecific CNV occurred across
401 each of the three genomes sequenced from hybridizable species (Pinosio *et al. 2016*).
402 This and similar studies suggest that a single genome assembly from a single locality is
403 likely not representative of the copy number diversity that exists within the sampled
404 population (Hirsch et al. 2014; Golicz et al. 2016; Gordon et al. 2017; Zhao et al. 2018).
405
406 By this notion, the following observations are in no way confirmatory without
407 additional sources of evidence. Although this dataset does not resolve interspecific
408 diversity, it is still representative of the diversity in Juglans, and allows for exceptional
409 resolution of patterns in gene family evolution. Multiple samples within sect.
410 *Rhysocaryon* and sect. *Dioscaryon*, and careful attention to informatic strategies, increases
411 confidence that the observed patterns across these genomes are true and not artifacts.
412
413 *Disease Resistance*
414
415 <u>*Losses in Dioscaryon*</u>
416 The absence of *Dioscaryon* gene models in the reticuline oxidase (berberine bridge
417 enzyme, BBE) annotated orthogroups shows a contraction before their divergence 22
418 MYA (Stevens *et al.* 2018). Enzymes in this family have been shown to contribute to
419 alkaloid production (Fujii *et al.* 2007) in California poppy (*Eschscholzia californica*) and
420 have been implicated in monolignol metabolism. Extreme (400-fold) upregulation of
421 enzymes in this family has been observed during pathogen attack and osmotic stress in
422 *Arabidopsis* (Daniel *et al.* 2015). Recent work in *Arabidopsis* demonstrated the function of

10

423    one BBE-like enzyme in oxidizing oligogalacturonides (OGs) and thereby diminishing
424    their elicitor activity (Benedetti *et al.* 2018). It is likely that the loss of the BBE gene
425    family in *J. regia* and *J. sigillata* occurred in the *Dioscaryon* ancestor but that does not
426    eliminate the possibility that these species were favored and therefore selected for their
427    potentially tamed secondary metabolite profiles. Until recently, chemical analyses in
428    *Juglans* have been limited to observational studies and comparisons of different
429    cultivars within a species (Vu *et al.* 2018; Vu *et al.* 2019). Additional studies contrasting
430    metabolomic profiles of domesticated species with their wild relatives will offer
431    valuable insight into tree domestication, especially when paired with genome
432    annotations.
433
434    The wall-associated kinases (WAKs) are a family of transmembrane receptor-like proteins
435    that bind pectin in the extracellular matrix (ECM) (Wagner and Kohorn, 2001). They are
436    necessary for cell expansion in Arabidopsis seedlings, but when bound to OGs also function
437    in defense response through Enhanced disease susceptibility 1 (EDS1) and Phytoalexin
438    deficient 4 (PAD4) dependent activation of MPK6-dependent pathway (Kohorn et al. 2009;
439    Brutus et al. 2010; Kohorn et al. 2014; Davidsson et al. 2017). Recent studies of WAKs have
440    shed light on their role in plant response to abiotic stressors (Marakli and Gozukirmizi.
441    2018, Xia et al. 2018) but many WAK family genes remain without functional
442    characterization. Because of this, the parallel contraction and loss of *J. hindsii* and *J. sigillata*
443    genes from multiple WAK annotated orthogroups is difficult to speculate on. A more
444    elaborate depiction of the WAK gene family will certainly shed light on the significance of
445    these losses. It is interesting to note, however, that two gene families (WAK and BBE) which
446    have members known to interact with OGs are both contracted in J. sigillata. These losses
447    suggest a significant shift in *J. sigillata* effector-triggered immunity.
448
449    _Losses and contractions in J. hindsii_
450    In California, the cultivation of *J. regia* is most commonly facilitated using Paradox
451    rootstock (*J. hindsii* ♀ × *J. regia* ♂), which is valued for its resistance to soil-borne
452    pathogens (Browne *et al.* 2015, Potter *et al.* 2002). Despite higher resistance to several
453    diseases, Paradox rootstock remain susceptible to *Armillaria* root rot, which is caused by
454    a basidiomycete, *A. mellea* in California (Baumgartner *et al.* 2013). The impact of this
455    disease is worsened by the lack of post-infection controls. Accordingly, discovering
456    resistant Paradox hybrids has been the focus of some research, but has achieved limited
457    success relative to the levels of *Armillaria* resistance reported in *J. hindsii* (Drakulic *et al.*
458    2017). Full genome annotations for these *Juglans* and others might be able to impart
459    clues about the genetic distinctions that contribute to these agriculturally interesting
460    phenotypes.
461

462    For comparison, interactions between *Arabidopsis* and the fungal pathogen, *Fusarium*
463    *oxysporum* are intensely studied as a model for plant fungal diseases. In this system, *F.*
464    *oxysporum* infections are potentiated by the alkalinization of soil around host root tissue
465    caused by plant RALF-triggered alkalinization response to pathogen secreted peptides,
466    RALFs (Rapid Alkalinization Factors) homologs (Masachis *et al.* 2016). These fungal
467    peptides target various members of transmembrane receptor-like kinases encoded by
468    the plant *Catharanthus roseus* Receptor-like Kinase (CrRLK1L) gene family. There are 17
469    reported CrRLK1L protein orthologs in *Arabidopsis* that have been implicated in a
470    variety of processes including immunity signaling, abiotic stress response and cell wall
471    dynamics (Kessler *et al.* 2010, Richter *et al.* 2017, Guo *et al.* 2018, Richter *et al.* 2018).
472    Several Basidiomycete genomes have been reported to encode RALF homologs, making
473    it plausible that *Armillaria* is among the fungi that utilize RALF-homolog effectors in
474    infection.
475
476    The current literature suggests a central role for CrRLK1L with respect to *F. oxysporum*
477    resistance. It is possible that the reduction or absence of the CrRLK1L orthologs
478    (annotated as FERONIA) in *J. hindsii* is at least partially responsible for its resistance to
479    *A. mellea* infection. The absence of *J. hindsii* models across five orthogroups annotated as
480    receptor-like protein kinase FERONIA-like in the Juglandaceae comparison
481    (OG0000687, OG0000471, OG0022392, OG0009045, OG0013911) including two for which
482    every other species is represented (OG0000687, OG0000471) warrants further
483    investigation. If the gene family is lost in *J. hindsii*, discovering any compensatory
484    mechanisms that might maintain the integrity of CrRLK1L-involved pathways could
485    have application in engineering fungus-resistant plants.
486
487    Like the observation that *Arabidopsis* FERONIA knockouts are more resistant to
488    *Fusarium* infection (Masachis *et al.* 2016), the loss of function *Arabidopsis* mutations in
489    RWA2 (reduced wall acetylation-2) led to increased resistance against the Ascomycete
490    pathogen, *Botrytis cinerea*, the causal agent of grey mold (Manabe *et al.* 2011). *B. cinerea*
491    belongs to a family of fungi, Botryosphaeriaceae, several of which are known to infect
492    the nuts of *J. regia* and related species (Moral et al. 2010). The Juglandaceae RWA2
493    orthogroup (OG0000145) was missing *J. hindsii* gene models. RWA2 is involved in
494    secondary cell wall synthesis and is regulated by SND1 (secondary wall-associated
495    NAC domain protein 1) (Lee *et al.* 2011). No experiments to date have assessed the
496    susceptibility of various *Juglans* species to *B. cinerea*, but it would be interesting to
497    examine resistance to the pathogen in a species without RWA2. The observed loss of
498    Chitinase-3 (Cht3) (OG0001222) in *J. hindsii* is consistent with the loss of FERONIA and
499    RWA2. Cht3 (along with glucanase and thaumatin-like protein) are aspects of plant
500    response to fungal invasion (Singh *et al.* 2012) and was found to be under positive
501    selection in the additional *Juglans* species (Table S7). If the loss of FERONIA and RWA2

502   do correspond to a weakened compatible host signature, the decreased incidence of
503   fungal infection would render such defense responses inessential.
504
505   *Mutation rate estimates and evidence of WGD*
506   Testing for WGD supported the hypothesis of a Juglandoid duplication. We observed
507   similar bimodal distributions of Ks values among syntenic blocks of genes in each of the
508   *Juglans-Pterocarya* pairwise analyses (Figure 4A). The bimodal distribution can be
509   attributed to a mixture of estimates from two distinct lineages; comparisons between
510   orthologous genes; and comparisons between more distant paralogous genes arising
511   from the whole genome duplication. Bimodal distributions for all *Juglans-Pterocarya*
512   pairwise comparisons are consistent with the WGD occurring prior to the radiation of
513   Juglans (Luo et al. 2015; Zhu et al. 2019). As additional confirmation, the most
514   prominent feature in the same analysis against the annotated genome of *Quercus robur*
515   is a peak at divergence values greater than those estimated for the WGD (Table S5).
516
517   Using the larger mode for each of the five distributions, we can estimate the nucleotide
518   substitution rate using the method of Zhu *et al.* (2019), for comparison. Using 66 MYA
519   as the assumed date of the WGD from Zhu *et al.* (2019), we obtain a synonymous
520   mutation rate of 2.7x10-9. This rate is higher than the rate of 2.3x10-9 estimated using 14
521   genes in Luo et al. (2015) and closer to the more recent estimate of 2.5x10-9 in Zhu *et al.*
522   (2019) using thousands of genes in a *J.regia* x *J.microcarpa* hybrid.  Our faster rate is still
523   more consistent with the rates of other woody perennials (e.g. Palm (Gaut *et al.* 1996)
524   and *P. trichocarpa* (Tuskan *et al.* 2006)), and still five times slower than the rate reported
525   for *Arabidopsis* (Koch *et al.* 2000).
526
527   A surprising observation was the distance between the two modes. We assume that the
528   estimated Ks of the smaller mode represents the between species divergence. The ratios
529   of the larger to smaller modes ranged from 6.5 to 7.3.  Interpretation of fossil data
530   (Manchester 1987) placed the initial split into *Rhysocaryon* and *Cardiocaryon* around 45
531   MYA, resolving around 38 MYA. Much closer to the assumed time of the WGD than
532   our bimodal distributions of Ks would indicate under a molecular clock. The
533   discrepancy in estimated WGD times may be due to the non-neutral nature of these
534   substitutions and departure from a molecular clock. However, a relevant observation
535   was recently made using a coalescent based approach. Bai *et al.* (2018) noted that
536   convergence of effective population size indicates a much earlier beginning for the
537   divergence among Juglans lineages. Our data could also be interpreted to support a
538   more recent divergence of walnut lineages.
539
540   The resources and services provided by *Juglans* species are nutritionally and culturally
541   significant. Their wood, used to construct furnishing and musical instruments, is valued

542 among woodworkers. Ink from walnut husks was used by Leonardo da Vinci and

543 Rembrandt. Brown dye from walnut stained fabrics was used in classical Rome,

544 medieval Europe, Byzantium and the Ottoman Empire. The genus is elevated in poetry

545 across the globe, including for its non-monetary benefits in Mary Oliver's "The Black

546 Walnut Tree" (Oliver, 1992) and nutritional properties in Tatsuji Miyoshi's "In Praise of

547 a Walnut" (Miyoshi, 1946). We are enthusiastic to contribute to the understanding of

548 and appreciation for this genus by constructing these genome annotation resources.

549

550 **METHODS**

551 *Repeat Library Generation and Softmasking*

552 The seven assemblies, ranging in size from 600 Mb to just under 1 Gb (2n=32) were

553 assessed for repeat content (Stevens *et al.* 2018).  Scaffolds and contigs less than 3Kbp in

554 length were removed from the assemblies prior to annotation.  RepeatModeler (v1.0.8)

555 was used to construct a repeat library through a combination of *de novo* and structural

556 prediction tools wrapped into the pipeline (Smit and Hubley, 2008). RepeatModeler

557 provided base annotations for the repeat elements (Table S1) and generated a consensus

558 library that was used as input to Repeatmasker (v4.0.6) to generate softmasked

559 genomes (Smit *et al.* 2013).

560

561 *Structural Annotation*

562 After softmasking, a set of 19 independent *J. regia* tissue-specific libraries described in

563 Chakraborty *et al* (2015) were aligned to the reference genomes via TopHat2 (v2.1.1)

564 (Kim *et al.* 2013).  The Illumina 85bp PE sequences were independently quality

565 controlled for a minimum length of 45bp and a minimum Phred-scaled quality score of

566 35 via Sickle (v. 1.33) prior to alignment.  Independent alignment files were sorted and

567 provided to Braker2 (v2.0) which generated a hints file for semi-supervised training of

568 the *ab initio* gene prediction package, Augustus (Stanke *et al.* 2008).  Braker2 utilizes

569 RNA-Seq reads directly to inform gene prediction and deduce the final models (Hoff *et*

570 *a*l. 2016).  The annotation files (GFF) produced were processed by gFACs, to filter out

571 incomplete or improbable gene models on the basis of completion (identifiable start and

572 stop codons) and canonical gene structure (micro-exons and micro-introns < 20bp are

573 filtered to reduce erroneous models).  The gFACs package also resolves conflicting

574 models and reports splice site statistics as well as other basic gene structure statistics

575 (Caballero and Wegrzyn, 2019).

576

577 *Functional Annotation*

578 The EnTAP functional annotation package was employed to remove unlikely gene

579 models and provide provisional functional information (Hart *et al.* 2019). Multi-exonic

580 and mono-exonic gene models were subjected to different functional filtering pipelines

581 that each utilized EnTAP. For multi-exonic genes, EnTAP (v 0.8.1) was provided three

582    curated databases (NCBI's Plant Protein (release 87), NCBI's RefSeq Protein (release 87),
583    and UniProtKB/Swiss-Prot) for similarity search (50% target and query coverage;
584    Diamond E-value .00001), followed by gene family assignment via the EggNOG
585    database and EggNOG-mapper toolbox (Jensen *et al.* 2008). Associations to gene
586    families provided the basis for Gene Ontology term assignment, identification of
587    protein domains (PFAM), and associated pathways (KEGG) (Finn *et al.* 2014; Ashburner
588    *et al.* 2017). Multi-exonics were removed from the set if they had neither sequence
589    similarity search result nor gene family assignment. Mono-exonic genes are typically
590    over-estimated in the process of *ab initio* genome annotation. To reduce this effect, they
591    were aligned to a custom curated database of monoexonic genes from other plant
592    species using 80% query coverage and 80% target coverage cutoffs in an independent
593    similarity search through EnTAP. EggNOG and PFAM were used in mono-exonic gene
594    model filtering as they were for multi-exonic filtering. After the first round of filters,
595    InterProScan (v5.25) was used to confirm gene family assignment and protein domains
596    in monoexonic gene models. Gene models without InterProScan annotations were
597    removed from the monoexonic set. For each species, the filtered multi-exonic and
598    mono-exonic gene sets were combined and passed back to gFACs to generate a
599    statistical profile and consistent annotation file in gene transfer format (GTF). Finally,
600    gene models that annotated with domains specific to retroelements were further filtered
601    from the final annotations based upon Pfam database descriptions. The entire set of
602    filtered gene models was evaluated for completeness. BUSCO (v3.0.2) was used with
603    default parameters and the embryophyta reference set of 1440 orthologs for this
604    purpose (Simão *et al.* 2015). Using the output from Augustus, we used gFACs to also
605    capture partial gene models. These were also functionally annotated used EnTAP, and
606    then compared using the same BUSCO analysis and ortholog set.
607
608    *Gene Family Classification and Evolution*
609    The proteins derived from the filtered genome annotations of each species were
610    processed with OrthoFinder-Diamond (v1.1.10) to provide information about
611    orthologous gene families. OrthoFinder is robust to incomplete models, differing gene
612    lengths, and larger phylogenetic distances (Emms and Kelly, 2015). Gene families
613    (orthogroups) in OrthoFinder are defined as homologous genes descended from a
614    single gene from the last common ancestor of the species examined. It is assumed that a
615    parental gene of each orthogroup was present in the common ancestor of the seven
616    species investigated. Two independent runs were conducted with OrthoFinder: *Juglans*
617    with the *Pterocarya* outgroup, and another that included these species with a set of 6
618    selected Eurosids (*Citrus grandis*, *Eucalyptus grandis*, *Arabidopsis thaliana*, *Carica papaya*,
619    *Populus trichocarpa* and *Quercus robur*). Rates of gene family evolution were calculated
620    for each orthogroup using the stochastic birth and death rate modeling implemented in
621    CAFE (v4.1) (De Bie *et al.* 2006). Species trees were constructed by applying estimated

15

622  divergence times from literature detailing rosid phylogeny to the known topology
623  (Magallón *et al.* 2014, Dong *et al.* 2017). Large variance in gene copy number between
624  species can lead to inaccurate calculation of birth and death rate parameters, therefore
625  large orthogroups with more than 100 gene models were removed prior to the analyses
626  and later analyzed separately using those parameters calculated by including only
627  orthogroups with < 100 gene models. Orthogroups represented by a single set of
628  paralogs were also removed because they are uninformative. Rapidly evolving gene
629  families (orthogroups) were identified using CAFE, which models the rate of gene
630  family evolution while accounting for the uncertainty in membership that results from
631  imperfect genome annotation. For each set, the lambda (birth and death rate) parameter
632  was calculated uniformly across the phylogeny. Orthogroups with a large size variance
633  among taxa were selected using a CAFE family-wide P-values <0.05. Those orthogroups
634  with accelerated rates of evolution were selected using branch-specific Viterbi P-values
635  <0.05. The gene-family losses described were independently confirmed using Exonerate
636  protein2genome alignments of the longest gene in the orthogroup to the genome of the
637  excluded species (90% similarity and score 1000) (Slater and Birney, 2005).
638  
639  Functional enrichment of rapidly evolving gene families was assessed independently
640  for each node and leaf of the Juglandaceae cladogram and across the entire set.
641  EggNOG gene descriptions of the longest gene model from each orthogroup were
642  compiled into a functional background. The gene model annotations from sets of
643  orthogroups found to be either rapidly expanding or rapidly contracting at each leaf or
644  node were compared to that background to estimate functional enrichment within the
645  set.
646  
647  *Selection Analysis*
648  To test for positive selection in gene families of interest, the coding sequence of gene
649  models from each orthogroup were iteratively clustered with USEARCH (v 9.0.2132) at
650  various identities beginning at 0.95 down to a minimum of 0.7 at intervals of 0.05.
651  Iterative clustering was terminated once a cluster with sufficient species representation
652  (relative to the species representation of that particular orthogroup) was produced and
653  chosen for use in selection analysis. A multiple sequence alignment of the longest gene
654  model from each species in that cluster was produced using Clustal Omega (v 1.2.4).
655  The multiple sequence alignments and species tree were provided to CODEML from
656  PAML (v 4.9) to calculate $\omega$ (dN/dS), the ratio of non-synonymous to synonymous
657  amino acid substitutions, across two models of adaptive evolution, including nearly
658  neutral and positive selection and the corresponding likelihood values. A likelihood
659  ratio test was used to determine the best model for each orthogroup.
660  
661  *Syntelog Analysis*

16

662   Genome alignment and analysis of syntenic genes was performed for each *Juglans*
663   genome against *Pterocarya stenoptera* using a CoGE (Lyons and Freeling, 2008)
664   SynMAP2 analysis. Genome alignment was performed using Last.  Five genes were
665   used as the minimum number of aligned pairs for DAGchainer (Haas *et al.* 2004).
666   Synonymous (Ks) and non-synonymous (Kn) coding sequence divergence was
667   estimated for syntenic protein coding gene pairs with CodeML (Yang, 2007).
668
669   **Data Availability:** The genomic resources described here are available at NCBI under
670   BioProject PRJNA445704 and the transcriptomic resources under BioProject
671   PRJNA232394.  These resources are also accessible from hardwoodgenomics.org and
672   treegenesdb.org. Functional annotations, gene models and gene transfer format (gtf)
673   files are also available on treegenesdb.org. Scripts and detailed processes used for this
674   study are accessible on https://gitlab.com/tree-genome-annotation/Walnut_Annotation.
675

682

683   **Author Contributions:** DBN, CHL, AD and KAS envisioned the resource and generated
684   the sequence data.  JLW, KAS, SZ, AJT and TF outlined the comparative methodology.
685   AJT, TF, SZ, MC and KAS analyzed the data. AJT, TF and JLW wrote the paper.

686    **Ashburner, M., Ball, C.A., Blake, J.A., et al.** (2000) Gene ontology: tool for the
687    unification of biology. The Gene Ontology Consortium. *Nat. Genet.,* **25**, 25–29.

688    **Bai, W.-N., Yan, P.-C., Zhang, B.-W., Woeste, K.E., Lin, K. and Zhang, D.-Y.** (2018)
689    Demographically idiosyncratic responses to climate change and rapid Pleistocene
690    diversification of the walnut genus *Juglans* (Juglandaceae) revealed by whole-genome
691    sequences. *New Phytologist,* **217**, 1726–1736.

692    **Baumgartner, K., Fujiyoshi, P., Browne, G.T., Leslie, C. and Kluepfel, D.A.** (2013)
693    Evaluating Paradox Walnut Rootstocks for Resistance to Armillaria Root Disease.
694    *HortScience,* **48**, 68–72.

695    **Beineke, W.F.** (1983) The Genetic Improvement of Black Walnut for Timber Production.
696    In J. Janick, ed. *Plant Breeding Reviews: Volume 1.* Boston, MA: Springer US, pp. 236–266.

697    **Benedetti, M., Verrascina, I., Pontiggia, D., Locci, F., Mattei, B., Lorenzo, G.D. and**
698    **Cervone, F.** (2018) Four Arabidopsis berberine bridge enzyme-like proteins are specific
699    oxidases that inactivate the elicitor-active oligogalacturonides. *The Plant Journal,* **94**, 260–
700    273.

701    **Bernard, A., Lheureux, F. and Dirlewanger, E.** (2017) Walnut: past and future of genetic
702    improvement. *Tree Genetics & Genomes,* **14**, 1.

703    **Browne, G.T., Grant, J.A., Schmidt, L.S., Leslie, C.A. and McGranahan, G.H.** (2011)
704    Resistance to Phytophthora and Graft Compatibility with Persian Walnut among
705    Selections of Chinese Wingnut. *HortScience,* **46**, 371–376.

706    **Browne, G.T., Leslie, C.A., Grant, J.A., Bhat, R.G., Schmidt, L.S., Hackett, W.P.,**
707    **Kluepfel, D.A., Robinson, R. and McGranahan, G.H.** (2015) Resistance to Species of
708    Phytophthora Identified among Clones of *Juglans microcarpa × J. regia. HortScience,* **50**,
709    1136–1142.

710    **Brutus, A., Sicilia, F., Macone, A., Cervone, F. and Lorenzo, G.D.** (2010) A domain
711    swap approach reveals a role of the plant wall-associated kinase 1 (WAK1) as a receptor
712    of oligogalacturonides. *PNAS,* **107**, 9452–9457.

713 **Buzo, T., McKenna, J., Kaku, S., Anwar, S.A. and McKenry, M.V.** (2009) VX211, A

714 Vigorous New Walnut Hybrid Clone with Nematode Tolerance and a Useful Resistance

715 Mechanism. *J Nematol*, **41**, 211–216.

716 **Caballero, M. and Wegrzyn, J.** (2019) gFACs: Gene Filtering, Analysis, and Conversion

717 to Unify Genome Annotations Across Alignment and Gene Prediction Frameworks.

718 *Genomics, Proteomics & Bioinformatics*.

719 **Callard, D., Axelos, M. and Mazzolini, L.** (1996) Novel Molecular Markers for Late

720 Phases of the Growth Cycle of Arabidopsis thaliana Cell-Suspension Cultures Are

721 Expressed during Organ Senescence. *Plant Physiology*, **112**, 705–715.

722 **Chakraborty, S., Britton, M., Martínez-García, P.J. and Dandekar, A.M.** (2016) Deep

723 RNA-Seq profile reveals biodiversity, plant–microbe interactions and a large family of

724 NBS-LRR resistance genes in walnut (*Juglans regia*) tissues. *AMB Express*, **6**, 12.

725 **Chakraborty, S., Britton, M., Wegrzyn, J., et al.** (2015) YeATS - a tool suite for

726 analyzing RNA-seq derived transcriptome identifies a highly transcribed putative

727 extensin in heartwood/sapwood transition zone in black walnut. *F1000Res*, **4**.

728 **Charrier, G., Bonhomme, M., Lacointe, A. and Améglio, T.** (2011) Are budburst dates,

729 dormancy and cold acclimation in walnut trees (*Juglans regia* L.) under mainly

730 genotypic or environmental control? *Int J Biometeorol*, **55**, 763–774.

731 **Chen, J., Mao, L., Lu, W., Ying, T. and Luo, Z.** (2016) Transcriptome profiling of

732 postharvest strawberry fruit in response to exogenous auxin and abscisic acid. *Planta*,

733 **243**, 183–197.

734 **Cheng, C.-Y., Krishnakumar, V., Chan, A.P., Thibaud-Nissen, F., Schobel, S. and**

735 **Town, C.D.** (2017) Araport11: a complete reannotation of the Arabidopsis thaliana

736 reference genome. *Plant J.*, **89**, 789–804.

737 **Chin, C.-S., Peluso, P., Sedlazeck, F.J., et al.** (2016) Phased diploid genome assembly

738 with single-molecule real-time sequencing. *Nature Methods*, **13**, 1050–1054.

739 **Cook, D.E., Lee, T.G., Guo, X., et al.** (2012) Copy Number Variation of Multiple Genes

740 at Rhg1 Mediates Nematode Resistance in Soybean. *Science*, **338**, 1206–1209.

741    **Cui, B., Pan, Q., Clarke, D., Villarreal, M.O., Umbreen, S., Yuan, B., Shan, W., Jiang, J.**

742    **and Loake, G.J.** (2018) S -nitrosylation of the zinc finger protein SRG1 regulates plant

743    immunity. *Nature Communications*, **9**, 4226.

744    **Daniel, B., Pavkov-Keller, T., Steiner, B., et al.** (2015) Oxidation of Monolignols by

745    Members of the Berberine Bridge Enzyme Family Suggests a Role in Plant Cell Wall

746    Metabolism. *J. Biol. Chem.*, **290**, 18770–18781.

747    **Davidsson, P., Broberg, M., Kariola, T., Sipari, N., Pirhonen, M. and Palva, E.T.** (2017)

748    Short oligogalacturonides induce pathogen resistance-associated gene expression in

749    Arabidopsis thaliana. *BMC Plant Biol*, **17**.

750    **De Bie, T., Cristianini, N., Demuth, J.P. and Hahn, M.W.** (2006) CAFE: a

751    computational tool for the study of gene family evolution. *Bioinformatics*, **22**, 1269–1271.

752    **Denton, J.F., Lugo-Martinez, J., Tucker, A.E., Schrider, D.R., Warren, W.C. and Hahn,**

753    **M.W.** (2014) Extensive Error in the Number of Genes Inferred from Draft Genome

754    Assemblies. *PLOS Computational Biology*, **10**, e1003998.

755    **Dong, W., Xu, C., Li, W., Xie, X., Lu, Y., Liu, Y., Jin, X. and Suo, Z.** (2017) Phylogenetic

756    Resolution in *Juglans* Based on Complete Chloroplast Genomes and Nuclear DNA

757    Sequences. *Front. Plant Sci.*, **8**.

758    **Drakulic, J., Gorton, C., Perez-Sierra, A., Clover, G. and Beal, L.** (2017) Associations

759    Between Armillaria Species and Host Plants in U.K. Gardens. *Plant Disease*, **101**, 1903–

760    1909.

761    **Eckert, A.J., Maloney, P.E., Vogler, D.R., Jensen, C.E., Mix, A.D. and Neale, D.B.**

762    (2015) Local adaptation at fine spatial scales: an example from sugar pine (Pinus

763    lambertiana, Pinaceae). *Tree Genetics & Genomes*, **11**, 42.

764    **Emms, D.M. and Kelly, S.** (2015) OrthoFinder: solving fundamental biases in whole

765    genome comparisons dramatically improves orthogroup inference accuracy. *Genome*

766    *Biology*, **16**, 157.

767    **Famula, R.A., Richards, J.H., Famula, T.R. and Neale, D.B.** (2018) Association genetics

768    of carbon isotope discrimination and leaf morphology in a breeding population of

769    *Juglans regia* L. *Tree Genetics & Genomes*, **15**, 6.

770    **Finn, R.D., Bateman, A., Clements, J., et al.** (2014) Pfam: the protein families database.

771    *Nucleic Acids Res*, **42**, D222–D230.

772    **Fujii, N., Inui, T., Iwasa, K., Morishige, T. and Sato, F.** (2007) Knockdown of berberine

773    bridge enzyme by RNAi accumulates (S)-reticuline and activates a silent pathway in

774    cultured California poppy cells. *Transgenic Res*, **16**, 363–375.

775    **Gaut, B.S., Morton, B.R., McCaig, B.C. and Clegg, M.T.** (1996) Substitution rate

776    comparisons between grasses and palms: synonymous rate differences at the nuclear

777    gene Adh parallel rate differences at the plastid gene rbcL. *PNAS*, **93**, 10274–10279.

778    **Golicz, A.A., Bayer, P.E., Barker, G.C., et al.** (2016) The pangenome of an

779    agronomically important crop plant *Brassica oleracea*. *Nature Communications*, **7**, 13390.

780    **Gordon, S.P., Contreras-Moreira, B., Woods, D.P., et al.** (2017) Extensive gene content

781    variation in the Brachypodium distachyon pan-genome correlates with population

782    structure. *Nat Commun*, **8**, 1–13.

783    **Gunn, B.F., Aradhya, M., Salick, J.M., Miller, A.J., Yongping, Y., Lin, L. and Xian, H.**

784    (2010) Genetic variation in walnuts (*Juglans regia* and *J. sigillata*◎; Juglandaceae): Species

785    distinctions, human impacts, and the conservation of agrobiodiversity in Yunnan,

786    China. *American Journal of Botany*, **97**, 660–671.

787    **Guo, H., Nolan, T.M., Song, G., Liu, S., Xie, Z., Chen, J., Schnable, P.S., Walley, J.W.**

788    **and Yin, Y.** (2018) FERONIA Receptor Kinase Contributes to Plant Immunity by

789    Suppressing Jasmonic Acid Signaling in Arabidopsis thaliana. *Current Biology*, **28**, 3316-

790    3324.e6.

791    **Haas, B.J., Delcher, A.L., Wortman, J.R. and Salzberg, S.L.** (2004) DAGchainer: a tool

792    for mining segmental genome duplications and synteny. *Bioinformatics*, **20**, 3643–3646.

793    **Hart, A.J., Ginzburg, S., Xu, M. (Sam), Fisher, C.R., Rahmatpour, N., Mitton, J.B.,**

794    **Paul, R. and Wegrzyn, J.L.** (2019) EnTAP: Bringing Faster and Smarter Functional

795    Annotation to Non-Model Eukaryotic Transcriptomes. *BioRxiv*, 307868.

796    **Hastings, P., Lupski, J.R., Rosenberg, S.M. and Ira, G.** (2009) Mechanisms of change in

797    gene copy number. *Nat Rev Genet*, **10**, 551–564.

798    Hirsch, C.N., Foerster, J.M., Johnson, J.M., et al. (2014) Insights into the Maize Pan-
799    Genome and Pan-Transcriptome. *The Plant Cell*, **26**, 121–135.
800    Hoff, K.J., Lange, S., Lomsadze, A., Borodovsky, M. and Stanke, M. (2016) BRAKER1:
801    Unsupervised RNA-Seq-Based Genome Annotation with GeneMark-ET and
802    AUGUSTUS. *Bioinformatics*, **32**, 767–769.
803    Jensen, L.J., Julien, P., Kuhn, M., Mering, C. von, Muller, J., Doerks, T. and Bork, P.
804    (2008) eggNOG: automated construction and annotation of orthologous groups of
805    genes. *Nucleic Acids Res*, **36**, D250–D254.
806    Kessler, S.A., Shimosato-Asano, H., Keinath, N.F., Wuest, S.E., Ingram, G., Panstruga,
807    R. and Grossniklaus, U. (2010) Conserved Molecular Components for Pollen Tube
808    Reception and Fungal Invasion. *Science*, **330**, 968–971.
809    Kim, D., Pertea, G., Trapnell, C., Pimentel, H., Kelley, R. and Salzberg, S.L. (2013)
810    TopHat2: accurate alignment of transcriptomes in the presence of insertions, deletions
811    and gene fusions. *Genome Biol.*, **14**, R36.
812    Kim, S.-H., Lee, K.-S., Son, J.-K., Je, G.-H., Lee, J.-S., Lee, C.-H. and Cheong, C.-J.
813    (1998) Cytotoxic Compounds from the Roots of *Juglans mandshurica. J. Nat. Prod.*, **61**,
814    643–645.
815    Koch, M.A., Haubold, B. and Mitchell-Olds, T. (2000) Comparative Evolutionary
816    Analysis of Chalcone Synthase and Alcohol Dehydrogenase Loci in Arabidopsis,
817    Arabis, and Related Genera (Brassicaceae). *Mol Biol Evol*, **17**, 1483–1498.
818    Koh, S.-J., Choi, Y.-I., Kim, Y., Kim, Y.-S., Choi, S.W., Kim, J.W., Kim, B.G. and Lee,
819    K.L. (2018) Walnut phenolic extract inhibits nuclear factor kappaB signaling in intestinal
820    epithelial cells, and ameliorates experimental colitis and colitis-associated colon cancer
821    in mice. *Eur J Nutr*.
822    Kohorn, B.D., Johansen, S., Shishido, A., Todorova, T., Martinez, R., Defeo, E. and
823    Obregon, P. (2009) Pectin activation of MAP kinase and gene expression is WAK2
824    dependent. *The Plant Journal*, **60**, 974–982.
825    Kohorn, B.D., Kohorn, S.L., Saba, N.J. and Meco-Martinez, V. (2014) Requirement for
826    Pectin Methyl Esterase and Preference for Fragmented Over Native Pectins for Wall

827    Associated Kinase Activated, EDS1/PAD4 Dependent Stress Response in Arabidopsis. *J.*

828    *Biol. Chem.*, jbc.M114.567545.

829    **Lee, C., Teng, Q., Zhong, R. and Ye, Z.-H.** (2011) The Four Arabidopsis REDUCED

830    WALL ACETYLATION Genes are Expressed in Secondary Wall-Containing Cells and

831    Required for the Acetylation of Xylan. *Plant Cell Physiol*, **52**, 1289–1301.

832    **Luo, M.-C., You, F.M., Li, P., et al.** (2015) Synteny analysis in Rosids with a walnut

833    physical map reveals slow genome evolution in long-lived woody perennials. *BMC*

834    *Genomics*, **16**, 707.

835    **Lyons, E., Pedersen, B., Kane, J. and Freeling, M.** (2008) The Value of Nonmodel

836    Genomes and an Example Using SynMap Within CoGe to Dissect the Hexaploidy that

837    Predates the Rosids. *Tropical Plant Biol.*, **1**, 181–190.

838    **Magallón, S., Gómez-Acevedo, S., Sánchez-Reyes, L.L. and Hernández-Hernández, T.**

839    (2015) A metacalibrated time-tree documents the early rise of flowering plant

840    phylogenetic diversity. *New Phytologist*, **207**, 437–453.

841    **Manabe, Y., Nafisi, M., Verhertbruggen, Y., et al.** (2011) Loss-of-Function Mutation of

842    REDUCED WALL ACETYLATION2 in Arabidopsis Leads to Reduced Cell Wall

843    Acetylation and Increased Resistance to Botrytis cinerea. *Plant Physiology*, **155**, 1068–

844    1078.

845    **Manchester, S.R.** (1987) The fossil history of the Juglandaceae. *Monogr.Syst.Bot.Missouri*

846    *Bot.Gard.*, **21**, 1–137.

847    **Marakli, S. and Gozukirmizi, N.** (2018) Analyses of abiotic stress and brassinosteroid-

848    related some genes in barley roots grown under salinity stress and HBR treatments:

849    Expression profiles and phylogeny. *Plant Biosystems - An International Journal Dealing*

850    *with all Aspects of Plant Biology*, **152**, 324–332.

851    **Marrano, A., Martínez-García, P.J., Bianco, L., et al.** (2018) A new genomic tool for

852    walnut (*Juglans regia* L.): development and validation of the high-density Axiom$^{TM}$ *J.*

853    *regia* 700K SNP genotyping array. *Plant Biotechnology Journal*, **0**.

854    **Martínez-García, P.J., Famula, R.A., Leslie, C., McGranahan, G.H., Famula, T.R. and**

855    **Neale, D.B.** (2017) Predicting breeding values and genetic components using

856    generalized linear mixed models for categorical and continuous traits in walnut (*Juglans*

857    *regia*). *Tree Genetics & Genomes*, **13**, 109.

858    **Masachis, S., Segorbe, D., Turrà, D., et al.** (2016) A fungal pathogen secretes plant

859    alkalinizing peptides to increase infection. *Nature Microbiology*, **1**, 16043.

860    **McGranahan, G.H., Leslie, C.A., Uratsu, S.L., Martin, L.A. and Dandekar, A.M.** (1988)

861    Agrobacterium-Mediated Transformation of Walnut Somatic Embryos and

862    Regeneration of Transgenic Plants. *Bio/Technology*, **6**, 800.

863    **Ming, R., Hou, S., Feng, Y., et al.** (2008) The draft genome of the transgenic tropical

864    fruit tree papaya (Carica papaya Linnaeus). *Nature*, **452**, 991–996.

865    **Miyoshi, T.** In Praise of the Walnut.

866    http://www.poetryinternational.org/pi/site/poem/item/16382

867    **Moral, J., Muñoz-Díez, C., González, N., Trapero, A. and Michailides, T.J.** (2010)

868    Characterization and Pathogenicity of Botryosphaeriaceae Species Collected from Olive

869    and Other Hosts in Spain and California. *Phytopathology*, **100**, 1340–1351.

870    **Oliver, M.** (2004) *New and Selected Poems, Volume One* Reprint edition., Beacon Press.

871    **Patterson, E.L., Pettinga, D.J., Ravet, K., Neve, P. and Gaines, T.A.** (2018) Glyphosate

872    Resistance and EPSPS Gene Duplication: Convergent Evolution in Multiple Plant

873    Species. *J Hered*, **109**, 117–125.

874    **Pina-Martins, F., Baptista, J., Pappas, G. and Paulo, O.S.** (2019) New insights into

875    adaptation and population structure of cork oak using genotyping by sequencing.

876    *Global Change Biology*, **25**, 337–350.

877    **Plomion, C., Aury, J.-M., Amselem, J., et al.** (2018) Oak genome reveals facets of long

878    lifespan. *Nature Plants*, **4**, 440.

879    **Pinosio, S., Giacomello, S., Faivre-Rampant, P., et al.** (2016) Characterization of the

880    Poplar Pan-Genome by Genome-Wide Identification of Structural Variation. *Mol Biol*

881    *Evol*, **33**, 2706–2719.

882    **Pollegioni, P., Woeste, K.E., Chiocchini, F., Olimpieri, I., Tortolano, V., Clark, J.,**

883    **Hemery, G.E., Mapelli, S. and Malvolti, M.E.** (2014) Landscape genetics of Persian

884    walnut (*Juglans regia* L.) across its Asian range. *Tree Genetics & Genomes*, **10**, 1027–1043.

885  **Pope, K.S., Dose, V., Silva, D.D., Brown, P.H., Leslie, C.A. and DeJong, T.M.** (2013)

886  Detecting nonlinear response of spring phenology to climate change by Bayesian

887  analysis. *Global Change Biology*, **19**, 1518–1525.

888  **Potter, D., Gao, F., Baggett, S., McKenna, J.R. and McGranahan, G.H.** (2002) Defining

889  the sources of Paradox: DNA sequence markers for North American walnut (*Juglans* L.)

890  species and hybrids. *Scientia Horticulturae*, **94**, 157–170.

891  **Prunier, J., Giguère, I., Ryan, N., Guy, R., Soolanayakanahally, R., Isabel, N., MacKay,**

892  **J. and Porth, I.** (2019) Gene copy number variations involved in balsam poplar (Populus

893  balsamifera L.) adaptive variations. *Molecular Ecology*, **28**, 1476–1490.

894  **Richter, J., Ploderer, M., Mongelard, G., Gutierrez, L. and Hauser, M.-T.** (2017) Role of

895  CrRLK1L Cell Wall Sensors HERCULES1 and 2, THESEUS1, and FERONIA in Growth

896  Adaptation Triggered by Heavy Metals and Trace Elements. *Front Plant Sci*, **8**.

897  **Richter, J., Watson, J.M., Stasnik, P., Borowska, M., Neuhold, J., Berger, M., Stolt-**

898  **Bergner, P., Schoft, V. and Hauser, M.-T.** (2018) Multiplex mutagenesis of four

899  clustered CrRLK1L with CRISPR/Cas9 exposes their growth regulatory roles in

900  response to metal ions. *Scientific Reports*, **8**, 12182.

901  **Settle, J., M., S. and Gonso, C.** (2015) 2015 Indiana Forest Products Price Report and

902  Trend Analysis. *Purdue Univ.,Dept. For. Nat. Resour.*, **October**, 17.

903  **Simão, F.A., Waterhouse, R.M., Ioannidis, P., Kriventseva, E.V. and Zdobnov, E.M.**

904  (2015) BUSCO: assessing genome assembly and annotation completeness with single-

905  copy orthologs. *Bioinformatics*, **31**, 3210–3212.

906  **Singh, D., Bhaganagare, G., Bandopadhyay, R., Prabhu, K.V., Gupta, P.K. and**

907  **Mukhopadhyay, K.** (2012) Targeted spatio-temporal expression based characterization

908  of state of infection and time-point of maximum defense in wheat NILs during leaf rust

909  infection. *Mol Biol Rep*, **39**, 9373–9382.

910  **Singh, R., Ong-Abdullah, M., Low, E.-T.L., et al.** (2013) Oil palm genome sequence

911  reveals divergence of interfertile species in Old and New worlds. *Nature*, **500**, 335–339.

912  **Slater, G.S.C. and Birney, E.** (2005) Automated generation of heuristics for biological

913  sequence comparison. *BMC Bioinformatics*, **6**, 31.

914    **Smit, A., Hubley, R. and Green, P.** RepeatMasker.

915    http://www.repeatmasker.org/

916    **Stanke, M., Diekhans, M., Baertsch, R. and Haussler, D.** (2008) Using native and

917    syntenically mapped cDNA alignments to improve de novo gene finding. *Bioinformatics*,

918    **24**, 637–644.

919    **Stevens, K.A., Woeste, K., Chakraborty, S., et al.** (2018) Genomic Variation Among and

920    Within Six *Juglans* Species. *G3: Genes, Genomes, Genetics*, **8**, 2153–2165.

921    **Sui, S., Luo, J., Liu, D., Ma, J., Men, W., Fan, L., Bai, Y. and Li, M.** (2015) Effects of

922    Hormone Treatments on Cut Flower Opening and Senescence in Wintersweet

923    (Chimonanthus praecox). *HortScience*, **50**, 1365–1369.

924    **Tuskan, G.A., Difazio, S., Jansson, S., et al.** (2006) The genome of black cottonwood,

925    Populus trichocarpa (Torr. & Gray). *Science*, **313**, 1596–1604.

926    **Tuskan, G.A., Groover, A.T., Schmutz, J., et al.** (2018) Hardwood Tree Genomics:

927    Unlocking Woody Plant Biology. *Front. Plant Sci.*, **9**.

928    **Van Bel, M., Bucchini, F. and Vandepoele, K.** (2019) Gene space completeness in

929    complex plant genomes. *Current Opinion in Plant Biology*, **48**, 9–17.

930    **Vijay, N., Poelstra, J.W., Künstner, A. and Wolf, J.B.W.** (2013) Challenges and

931    strategies in transcriptome assembly and differential gene expression quantification. A

932    comprehensive in silico assessment of RNA-seq experiments. *Molecular Ecology*, **22**, 620–

933    634.

934    **Vu, D.C., Lei, Z., Sumner, L.W., Coggeshall, M.V. and Lin, C.-H.** (2019) Identification

935    and quantification of phytosterols in black walnut kernels. *Journal of Food Composition*

936    *and Analysis*, **75**, 61–69.

937    **Vu, D.C., Vo, P.H., Coggeshall, M.V. and Lin, C.-H.** (2018) Identification and

938    Characterization of Phenolic Compounds in Black Walnut Kernels. *J. Agric. Food Chem.*,

939    **66**, 4503–4511.

940    **Wagner, T.A. and Kohorn, B.D.** (2001) Wall-Associated Kinases Are Expressed

941    throughout Plant Development and Are Required for Cell Expansion. *The Plant Cell*, **13**,

942    303–318.

9

943    **Wang, X., Xu, Y., Zhang, S., et al.** (2017) Genomic analyses of primitive, wild and
944    cultivated citrus provide insights into asexual reproduction. *Nature Genetics*, **49**, 765–
945    772.

946    **Weckerle, C., Huber, F.K., Yongping, Y. and Weibang, S.** (2005) Walnuts among the
947    Shuhi in Shuiluo, Eastern Himalayas. *Economic Botany*, **59**, 287–290.

948    **Würschum, T., Langer, S.M., Longin, C.F.H., Tucker, M.R. and Leiser, W.L.** (2018) A
949    three-component system incorporating Ppd-D1, copy number variation at Ppd-B1, and
950    numerous small-effect quantitative trait loci facilitates adaptation of heading time in
951    winter wheat cultivars of worldwide origin. *Plant, Cell & Environment*, **41**, 1407–1416.

952    **Xia, Y., Yin, S., Zhang, K., Shi, X., Lian, C., Zhang, H., Hu, Z. and Shen, Z.** (2018)
953    OsWAK11, a rice wall-associated kinase, regulates Cu detoxification by alteration the
954    immobilization of Cu in cell walls. *Environmental and Experimental Botany*, **150**, 99–105.

955    **Xu, H., Yu, X., Qu, S. and Sui, D.** (2013) Juglone, isolated from *Juglans mandshurica*
956    Maxim, induces apoptosis via down-regulation of AR expression in human prostate
957    cancer LNCaP cells. *Bioorganic & Medicinal Chemistry Letters*, **23**, 3631–3634.

958    **Yang, Z.** (2007) PAML 4: Phylogenetic Analysis by Maximum Likelihood. *Mol Biol Evol*,
959    **24**, 1586–1591.

960    **Yao, Y., Zhang, Y.-W., Sun, L.-G., et al.** (2012) Juglanthraquinone C, a novel natural
961    compound derived from *Juglans mandshurica* Maxim, induces S phase arrest and
962    apoptosis in HepG2 cells. *Apoptosis*, **17**, 832–841.

963    **Zhang, W., Jiao, Z., Shang, T. and Yang, Y.** (2015) Demography and spectrum analysis
964    of *Juglans cathayensis* populations at different altitudes in the west Tianshan valley in
965    Xinjiang, China. *Ying Yong Sheng Tai Xue Bao*, **26**, 1091–1098.

966    **Zhao, P., Zhou, H.-J., Potter, D., et al.** (2018) Population genetics, phylogenomics and
967    hybrid speciation of *Juglans* in China determined from whole chloroplast genomes,
968    transcriptomes, and genotyping-by-sequencing (GBS). *Molecular Phylogenetics and*
969    *Evolution*, **126**, 250–265.

970    **Zhao, Q., Feng, Q., Lu, H., et al.** (2018) Pan-genome analysis highlights the extent of
971    genomic variation in cultivated and wild rice. *Nat Genet*, **50**, 278–284.

972    **Zhu, T., Wang, L., You, F.M., et al.** (2019) Sequencing a Juglans regia × J. microcarpa

973    hybrid yields high-quality genome assemblies of parental species. *Hortic Res*, **6**, 1–16.

974

975

976

977

978

979

980

981

982

983

984

985

986

987

988

989

990

991

992

993

994

995

996

997

998

999

1000

1001

1002

1003

1004

1005

1006

1007

1008

1009

1010

1011    **Figure Legends:**

1012

1013    Figure 1: Approximate ranges of each annotated species with shapes denoting the three
1014    sections of *Juglans* and the outgroup genus, *Pterocarya.*

1015

1016    Figure 2: As a measure of gene annotation completeness, the gene models from the
1017    seven annotations were compared to a set of 1,440 embryophyta putative single-copy
1018    orthologs using BUSCO (Benchmarking Universal Single-copy Orthologs) (Table S3).
1019    Orthologs were either found once in the gene models (Complete Single), multiple times
1020    (Complete Duplicated), partially (Fragmented), or were not found at all (Missing).

1021

1022    Figure 3 (A) Distribution of species membership across orthogroups. Tiling beneath the
1023    histogram indicates the species contributing gene models to each orthogroup in the set.
1024    Set size is displayed as height on histogram. The horizontal histogram indicates the
1025    number of orthogroups found in each species. Blue indicates data from groups
1026    composed of *Rhysocaryon* species while green bars show Eurasian species (*Dioscaryon*
1027    and *Cardiocaryon*). B) Cladogram with associated stacked histogram reflecting the
1028    number of genes belonging to orthogroups specific to the color-indicated groups.

1029

1030    Figure 4 Histograms of substitution rates for coding genes determined by SynMAP to
1031    be syntenic between *Juglans hindsii* and *Pterocarya stenoptera*. Two peaks are visible in
1032    both the non-synonymous (A) and synonymous (B) distributions. In both cases the
1033    highlighted righthand peak represents the older WGD. Table S8 summarizes the
1034    distributions for all annotated *Juglans* genomes described here against *P. stenoptera*.

1035

1036    Figure 5: Phylogenetic tree constructed from divergence times in literature displaying
1037    numbers of expanded (blue) and contracted (red) orthogroups per terminal branch
1038    discovered using OrthoFinder/CAFÉ with the 13 species Eurosid analysis. The number
1039    of significant (P-value <0.05, Viterbi P-value <0.05) expansions and contractions at each
1040    node and leaf are shown in parentheses.

1041

1042
1043
1044
1045
1046
1047
1048
1049

## Tables:

Table 1:Genome assembly statistics for seven Juglandaceae species

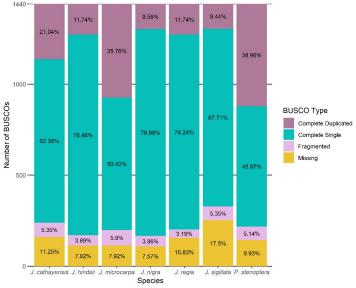|  | Juglans. hindsii | Juglans nigra | Juglans cathayensis | Juglans microcarpa | Juglans sigillata | Juglans regia | Pterocarya stenoptera |
|---|---|---|---|---|---|---|---|
| **Plant Info** | | | | | | | |
| **Name** | 'Rawlins' | 'Sparrow' | 'Wild Walnut' | '83-129' | 'Yangbi 1' | 'Chandler' | '83-13' |
| **Cultivar** | DJUG105 | A30 | DJUG11.03 | DJUG29.11 | DJUG951.04 | 64-172 | DPTE1.09 |
| **Source** | NCGR | MU | NCGR | NCGR | NCGR | UCD | NCGR |
| **Assembly** | | | | | | | |
| **Version** | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 | 1.1 | 1.0 |
| **Size (Mbp)** | 605.70 | 605.05 | 751.37 | 896.45 | 622.24 | 686.52 | 936.89 |
| **Scaffolds** | 273,094 | 232,579 | 332,634 | 329,873 | 282,224 | 186,636 | 396,056 |
| **N50 (Kbp)** | 512.79 | 118.45 | 158.25 | 145.01 | 218.35 | 278.29 | 159.70 |
| **Annotated Assembly (> 3Kbp scaffolds)** | | | | | | | |
| **Size (Mbp)** | 586.05 | 580.70 | 719.60 | 862.79 | 585.63 | 686.52 | 902.23 |
| **Scaffolds** | 4,672 | 5,896 | 10,342 | 12,024 | 6,413 | 11,848 | 11,574 |
| **N50 (Kbp)** | 540.03 | 271.37 | 168.53 | 151.91 | 238.36 | 278.30 | 167.10 |

Table 2: Structural and functional annotations for seven Juglandaceae species

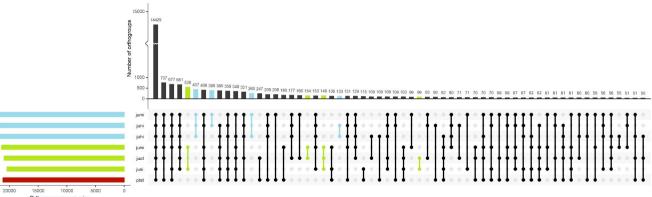|  | Juglans. hindsii | Juglans nigra | Juglans cathayensis | Juglans microcarpa | Juglans sigillata | Juglans regia | Pterocarya stenoptera |
|---|---|---|---|---|---|---|---|
| **Structural Annotation** | | | | | | | |
| Repeat Content | 46.97% | 47.34% | 48.03% | 46.87% | 46.69% | 47.96% | 43.89% |

13

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| Total Genes | 28,664 | 28,335 | 34,066 | 41,611 | 26,835 | 30,626 | 44,318 |
| Total Complete, Multi-exonics | 24,500 | 23,290 | 28,915 | 35,319 | 22,898 | 26,166 | 36,984 |
| Total Complete, mono-exons | 4,164 | 4,426 | 5,151 | 6,292 | 3,937 | 4,460 | 7,334 |
| Gene Length (Avg) | 4,406.38 | 4,301.45 | 4,193.80 | 4,008.01 | 4,373.40 | 4,235.02 | 3,944.81 |
| CDS Length (Avg) | 1,267.37 | 1,277.62 | 1,220.28 | 1,199.12 | 1,250.97 | 1,220.42 | 1203.44 |
| Exons per Gene (Average) | 6.30 | 6.38 | 6.06 | 5.98 | 6.32 | 6.14 | 5.91 |
| Canonical Splice Sites (%) | 98.70% | 98.67% | 98.69% | 98.70% | 98.60% | 98.77% | 98.76% |
| **Functional Annotation** | | | | | | | |
| EnTAP (Similarity Search) | 23,822 | 23,607 | 27,815 | 33,711 | 22,036 | 25,420 | 35,771 |
| EnTAP (Gene Family only) | 4,842 | 4,728 | 6,251 | 7,900 | 4,799 | 5,206 | 8.547 |

1055
1056

14

△ *Pterocarya*
□ Dioscaryon
○ Cardiocaryon
⬠ Rhysocaryon

*J. regia*

△ *P. stenoptera*

○ *J. cathayensis*

□ *J. sigillata*

⬠ *J. hindsii*

⬠ *J. nigra*

○ *J. microcarpa*

Number of genes

Species by *P. stenoptera*

— *J. cathayensis*
— *J. hindsii*
— *J. microcarpa*
— *J. nigra*
— *J. sigillata*

Gene Families
Expansion/Contraction
(significant)