1 **Whole genome duplication in coast redwood (*Sequoia sempervirens*) and its**

2 **implications for explaining the rarity of polyploidy in conifers**

3 Alison Dawn Scott, Noah Stenz, David A. Baum

4 Department of Botany, University of Wisconsin, Madison, 430 Lincoln Dr., Madison WI 53706

5 SUMMARY

6 • Whereas polyploidy is common and an important evolutionary factor in most land

7 plant lineages it is a real rarity in gymnosperms. Coast redwood (*Sequoia*

8 *sempervirens*) is the only hexaploid conifer and one of just two naturally

9 polyploid conifer species. Numerous hypotheses about the mechanism of

10 polyploidy in *Sequoia* and parental genome donors have been proffered over the

11 years, primarily based on morphological and cytological data, but it remains

12 unclear how *Sequoia* became polyploid and why this lineage overcame an

13 apparent gymnosperm barrier to whole-genome duplication (WGD).

14 • We sequenced transcriptomes and used phylogenetic inference, Bayesian

15 concordance analysis, and paralog age distributions to resolve relationships

16 among gene copies in hexaploid coast redwood and its close relatives.

17 • Our data show that hexaploidy in the coast redwood lineage is best explained by

18 autopolyploidy or, if there was allopolyploidy, this was restricted to within the

19 Californian redwood clade. We found that duplicate genes have more similar

20 sequences than would be expected given evidence from fossil guard cell size

21 which suggest that polyploidy dates to the Eocene.

22 • Conflict between molecular and fossil estimates of WGD can be explained if

23 diploidization occurred very slowly following whole genome duplication. We

24 extrapolate from this to suggest that the rarity of polyploidy in conifers may be

25 due to slow rates of diploidization in this clade.

26

27 KEYWORDS: whole genome duplication, polyploidy, *Sequoia sempervirens*, conifer,

28 gymnosperm

29

30

31

32   **INTRODUCTION**

33   Polyploidy has profound long- and short-term genetic consequences (Adams & Wendel,

34   2005; Otto & Whitton, 2000; etc.), and facilitates adaptive evolution (Soltis et al., 2008;

35   etc). Studies of genome sequences, expressed genes, and cytogenetics suggest that all

36   land plant lineages have experienced polyploidization in their evolutionary history,

37   though clades differ in the extent of recent whole genome duplication

38   (neopolyploidization). While there are thousands of neopolyploid mosses, ferns and

39   angiosperms, the phenomenon is relatively rare in gymnosperms, and especially conifers.

40   There are only two polyploid conifer species: alerce, *Fitzroya cupressoides* (4x), and

41   coast redwood, *Sequoia sempervirens* (6x). Why is polyploidy so rare in conifers? Does it

42   reflect rare formation of polyploid individuals, for example due to a lack of unreduced

43   gametes, or another barrier to allopolyploid formation? Or, do polyploid taxa form in

44   gymnosperms, but fail to give rise to successful clades? To shed light on these questions,

45   we studied the evolutionary history of coast redwood with the goal of determining when

46   polyploidy occurred and whether it entailed allopolyploidy.

47

48   Coast redwoods are long-lived trees (some over 2,000 years; Burns & Honkala, 1990)

49   that thrive in the foggy coastal forests of central and northern California. Coast redwoods

50   are among the world's tallest living trees (up to 115 meters; Ishii et al., 2014). *Sequoia* is

51   a monotypic genus whose closest relatives are the giant sequoia of the Californian Sierra

52   Nevada (*Sequoiadendron giganteum*) and the Chinese dawn redwood (*Metasequoia*

53   *glyptostroboides*). Though the three modern redwood species have distinct ranges, fossil

54   data suggest that diverse redwood lineages were widely distributed across the Northern

55   Hemisphere from the Cretaceous onwards (Miller, 1977). The oldest redwood fossils are

56   from South Manchuria (present-day China) and Boulogne-sur-Mer (northern France) and

57   date back to the mid-to-late Jurassic, suggesting the redwood clade is at least 146 million

58   years old (Zeiller and Fliche, 1903; Endo, 1951).

59   *Sequoidendron* and *Metasequoia* are diploids with 2n=22 (Schlarbaum and Tshuchiya,

60   1984). Hirayoshi and Nakamura (1943) first determined the correct chromosome number

61   of *Sequoia* and proved that it is a hexaploid with 2n=66. Hexaploidy in *Sequoia* was later

62   corroborated by Stebbins (1948), Saylor and Simons (1970) and Ahuja and Neale (2002).

63    Relying on the well-known correlation between guard cell size and genome size (e.g.,

64    Beaulieu et al., 2008), Miki and Hikita (1951) studied stomatal guard-cell size in Pliocene

65    fossils of Metasequoia and Sequoia. As fossil guard cells were the same size as extant

66    guard cells, Miki and Hikita concluded *Sequoia* has been hexaploid since at least the

67    Pliocene (2.5-5 million years ago). This estimate was pushed back significantly by Ma et

68    al. (2005), who describe fossils from the Eocene (33-53mya) with guard cells of a size

69    taken to indicate polyploidy.

70    Morphological similarities among modern redwoods led to hypotheses of allopolyploidy

71    in *Sequoia* involving hybridization between extinct diploid *Sequoia* and ancestors of

72    either *Metasequoia* (Stebbins, 1948) or *Sequoiadendron* (Doyle, 1945). Despite the

73    distance among their modern ranges, the overlap in fossil distributions of *Sequoia*,

74    *Sequoiadendron*, and *Metasequoia* make this hypothesis plausible. Another hypothesis is

75    that an extinct member of the Taxodiaceae, perhaps a member of *Taxodium*, contributed

76    to the hexaploid genome of *Sequoia* (Stebbins, 1948; Saylor and Simons, 1970). Ahuja

77    and Neale (2002), in contrast, suggested that the "missing" parent of *Sequoia* may have

78    been a member of the *Cryptomeria*, *Taiwania*, or *Athrotaxis* lineages.

79    Before the advent of molecular phylogenetics, auto- and allopolyploids were

80    distinguished by observing chromosome behavior during meiosis. Autopolyploidy

81    (generally interpreted as occurring within a single species) and allopolyploidy (involving

82    hybridization among species) represent extremes of a spectrum. Autopolyploids have

83    multiple sets of very similar homologous chromosomes, which tends to manifested

84    cytogenetically as the formation of multivalents (e.g. groups of four or six

85    chromosomes). Allopolyploids, in contrast, arise from the fusion of divergent genomes

86    which, in the extreme case results in bivalent formation by each homologous

87    chromosome, as observed in diploid organisms. However, chromosome pairing at

88    meiosis is rarely definitive as allopolyploidy can result in multivalent formation among

89    homeologs if hybridizing species are closely related, and bivalent formation is eventually

90    reestablished following autopolyploidy by the process of diploidization (Ramsey and

91    Schemske, 2002; Parisod et al., 2010).

92

93    In addition to cytogenetic lines of evidence, segregation patterns can be useful to

94    distinguish auto- and allopolyploids. An autopolyploid forming multivalents at meiosis

95    will produce equal frequencies of all possible allele combinations. In the case of *Sequoia*,

96    this pattern is called hexasomic inheritance. Allopolyploids do not typically form

97    multivalents at meiosis, resulting in simple disomic inheritance (as seen in diploids).

98    Again, these are only the most extreme possibilities, as both the diploidization process

99    and polyploidy involving a mixture of similar and divergent chromosomes (i.e. segmental

100   allopolyploidy sensu Stebbins) can lead to intermediate inheritance patterns.

101

102   Studies of meiotic chromosome pairing in *S. sempervirens* reported a mixture of bivalents

103   and multivalents (Stebbins, 1948; Schlarbaum and Tsuchiya, 1984; Ahuja and Neale,

104   2002). This led Stebbins (1948) and Schlarbaum and Tsuchiya (1984a, b) to suggest that

105   hexaploidy involved both auto- and allopolyploidy. A similar result was obtained by

106   Rogers (1997), who used allozymes to study inheritance patterns in *Sequoia*. However,

107   neither the pairing nor genetic data are sufficient to distinguish segmental allopolyploidy

108   from autoployploidy followed by partial diploidization  We set out to use modern

109   genomic approaches to revisit the evolutionary history of polyploidy in *S. sempervirens*

110   and see if, by doing so, we could also gain insights into why polyploidy is so rare in

111   gymnosperms.

112

113   **MATERIALS AND METHODS**

114   *Transcriptome sequencing and assembly*

115   Total RNA was extracted from foliage samples of *S. sempervirens*, *S. giganteum*, *M.*

116   *glyptostroboides*, and the outgroup *Thuja occidentalis* (eastern white cedar) with a

117   CTAB/Chisam extraction protocol followed by Qiagen RNeasy cleanup. Illumina TruSeq

118   cDNA libraries were prepared and sequenced on an Illumina HiSeq 2000 with 100bp

119   paired-end reads at either the UW Biotech Center (Madison, WI) or at the SciLife

120   Laboratory (Stockholm, Sweden).

121

122   *Sequence analysis and alignment*

123    We assembled raw reads *de novo* with Trinity vers. 2014-07-17 (Grabherr et al., 2011),

124    with default settings and Trimmomatic processing. After assembly, contigs were

125    translated using TransDecoder vers. 2014-07-04 (Haas et al., 2013;

126    http://transdecoder.sf.net) with a minimum protein length of 100aa. Translated contigs

127    were filtered using the Evigene pipeline vers. 2013.07.27

128    (http://arthropods.eugenes.org/EvidentialGene/about/EvidentialGene_trassembly_pip

129    e.html). Ortholog clusters shared among *S. sempervirens*, *S. giganteum*, *M.*

130    *glyptostroboides*, and *T. occidentalis* were identified using the translated transcriptome

131    assemblies by ProteinOrtho ver. 5.11 (Lechner et al., 2011), using an algebraic

132    connectivity cutoff of 0.25. Custom Perl scripts (available at github.com/nstenz) were

133    used to identify ortholog sets that contained a single copy in diploids (*S. giganteum*, *M.*

134    *glyptostroboides*, and *T. occidentalis*) and between one and three copies in the hexaploid

135    *S. sempervirens*. As these putatively single-copy protein-coding sequences show marked

136    conservation among species, we assumed that allelic variants would generally be

137    combined into a single contig. We used MUSCLE v. 3.8.13, 64bit  (Edgar, 2004a,b), with

138    default alignment settings to align the ortholog sets at the protein level before using a

139    custom PERL script to generate the corresponding nucleotide alignment.

140

141    *Single-variant gene trees and concordance analyses*

142    For each orthogroup that included only one sequence variant in *S. sempervirens* we

143    estimated phylogenetic trees using MrBayes vers. 3.2.2 64bit (Huelsenbeck & Ronquist,

144    2001; Ronquist & Huelsenbeck, 2003) with the settings: nst = 6; rates = invgamma; ngen

145    = 1.1 million; burnin = 100,000; samplefreq = 40; nruns = 4; nchains = 3; temp = 0.45;

146    swapfreq = 10. BUCKy vers. 1.4.4 (Ané et al., 2007; Larget et al., 2010) was then used to

147    estimate the proportion of genes that have each possible resolution in the redwood clade

148    while taking account of uncertainty in individual gene trees. Post-burnin posterior

149    distributions from MrBayes were combined in BUCKy for 1 million generations with $\alpha =$

150    1. All trees were rooted on the outgroup, *Thuja occidentalis.*

151

152    *Density distribution of $K_s$ estimates*

153   To build an age distribution of $K_s$ (the average number of synonymous substitutions per

154   synonymous site) within each transcriptome we identified duplicate genes using custom

155   Perl scripts (available at github.com/nstenz). Assembled contigs were translated using

156   TransDecoder with a minimum protein length of 100aa, as above. Duplicate genes were

157   identified using BLAT (Kent 2002) on translated contigs and then duplicate gene pairs

158   were aligned and back translated into their corresponding nucleotide sequence. We

159   estimated $K_s$ on each pair of nucleotide alignments using $K_aK_s$ calculator (model GY;

160   Zhang et al., 2006). We excluded $K_s$ values greater than 2 to avoid the effects of $K_s$

161   saturation, and plotted the resulting $K_s$ values in a density plot in R (R core team, 2013).

162   To identify significant features of the $K_s$ frequency distributions we used SiZer

163   (Chaudhuri and Marron, 1999).

164

165

166   *Multi-variant gene trees and tree-based $K_s$ estimates*

167   For alignments containing a single variant in diploid taxa and two or three variants in

168   hexaploid *Sequoia*, we estimated phylogenetic trees with raxml vers. 8.1.20 (100

169   bootstrap replicates; GTRGAMMA; Stamatakis, 2006). We then used PAML (Yang,

170   1997) to obtain a tree-based estimate of $K_s$. PAML calculates branch lengths along the

171   ML tree using a model that estimates the rate of synonymous and non-synonymous

172   substitutions ($D_s$ and $D_n$, respectively) separately for each branch. We imposed a

173   molecular clock assumption (clock=1) to obtain an ultrametric tree. By multiplying a

174   branch's length by its $D_s$ and summing over intervening branches between two tips we

175   could obtain an estimate of the patristic $K_s$ distance between *Sequoia* homeologs and how

176   this compares to the $K_s$ of copies from different species.

177

178   In order to obtain an approximate date for gene duplication, we divided the depth of the

179   gene duplication in $K_s$ units by an average mutation rate for conifers of $0.68 \times 10^{-9}$

180   synonymous substitutions per synonymous site per year (Buschiazzo et al., 2012).

181   *Sequoia* is hexaploid, so at least two whole genome duplications must have occurred in

182   the past. As each whole genome duplication event is expected to yield a normal

183   distribution of $K_s$ values, we used EMMIX v.1.3 (Mclachlan et al., 1999) to fit a mixture

184    model of normal distributions as a way to assign putative homeologs to each duplication

185    event and estimate their ages. We allowed EMMIX to fit 1-2 normal distributions, with

186    the optimal model selected based on AIC and BIC scores.

187

188    **RESULTS**

189    Our *de novo* transcriptome assemblies ranged from 70 to 101mbp in length (Table 1).

190    Assembled contigs per species ranged from 80,126 to 128,005.
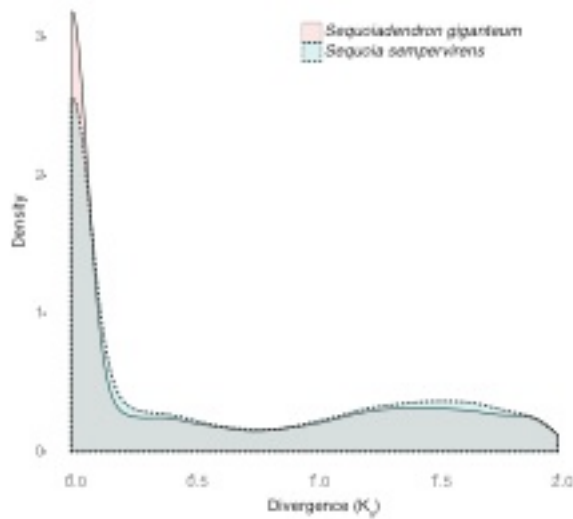
191

192    Table 1: Assembly statistics

| Taxon | Raw reads (paired end) | Assembly length (mbp) | Contigs | N50 |
|---|---|---|---|---|
| *Sequoia sempervirens* | 55,052,935 | 85.6 | 128,005 | 1,118 |
| *Sequoiadendron giganteum* | 56,665,524 | 101.3 | 115,519 | 1,619 |
| *Metasequoia glyptostroboides* | 29,502,075 | 78.6 | 83,120 | 1,668 |
| *Thuja occidentalis* | 31,116,702 | 70.0 | 80,126 | 1,607 |

193

194    Assuming synonymous substitutions happen at a constant rate over time, $K_s$ can be used

195    as a proxy for the age of duplicate genes. To estimate the distribution of pairwise $K_s$

196    distance within each genome, we identified all duplicate genes, which numbered 33,544,

197    39,236, and 26,485, in *S. sempervirens*, *S. giganteum*, and *M. glyptostroboides*,

198    respectively. Paralog age distribution plots for all three taxa revealed a peak at a $K_s \approx 1.5$,

199    of which those for *S. sempervirens*, *S. giganteum* are shown in Fig. 1. Allowing for the

200    approximate nature of these calculations, this peak likely corresponds to the seed plant

201    whole genome duplication previously dated at 319 Ma (Jiao et al., 2011).  Despite the

202    expectation that hexaploid *Sequoia* would have at least one other, much younger peak

203    corresponding to a polyploidization event in perhaps the Eocene (Ma et al., 2005), this

204    was not visible in the age distribution plots (Fig. 1). Results from SiZer also did not

205    indicate any significant peak unique to the *Sequoia* $K_s$ plot.

206

208

209    Figure 1: Density distribution of pairwise Ks between duplicate genes in *Sequoia* (pink) and

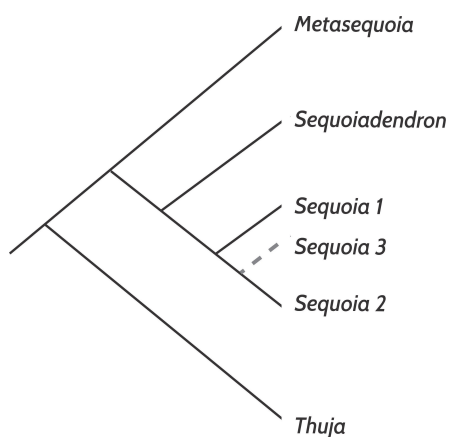210    *Sequoiadendron* (cyan).

211

212    To distinguish the evolutionary relationships among redwoods and look for evidence of

213    ancestral hybridization, we used Bayesian concordance analysis and estimated genomic

214    support for each of three possible topologies for an unrooted four-taxon tree. First we

215    built individual gene trees from 7,819 ortholog groups that each had one sequence variant

216    in each diploid species (*Sequioadendron, Metasequoia, Thuja*) and one, two, or three

217    sequence variants in the hexaploid, *Sequoia*. Alignment lengths in this set varied from

218    301-5,736 bp, with a median of 1,104. Of these alignments 7,602 included a single

219    *Sequoia* copy, whereas 217 included one or two *Sequoia* sequence variants. Among the

220    7,602 alignments that included a single copy in *S. sempervirens* the most frequently

221    supported topology placed *S. sempervirens* sister to *Sequoiadendron* (Fig. 2) with a

222    concordance factor (CF; Baum 2007) mean estimate of 0.79 and a 95% credibility

223    interval of 0.78-0.80. The two minor topologies (*Sequoia + Metasequoia*; *Metasequoia +*

224    *Sequoiadendron*) had concordance factors of 0.10(0.09-0.11) and 0.11(0.10, 0.12),

225    respectively (Fig. 2). These results show that, if *Sequoia* arose from allopolyploidy, it

226    only involved genome donors in the Californian redwood clade (i.e., the clade that

227    includes *S. sempervirens* and *Sequoiadendron*). However, autoploidy is also a possibility.

228

Figure 2: Bayesian concordance analysis of 7,602 gene trees. For each of three possible topologies, the concordance factor (proportion of loci in the sample having the clade) and its 95% credibility interval are shown.



Figure 3: Cladogram summarizing 184 gene trees as estimated by MrBayes.

In order to obtain estimates for the divergence of *Sequoia* duplicates relative to interspecies divergences and to re-evaluate evidence for allopolyploidy within the
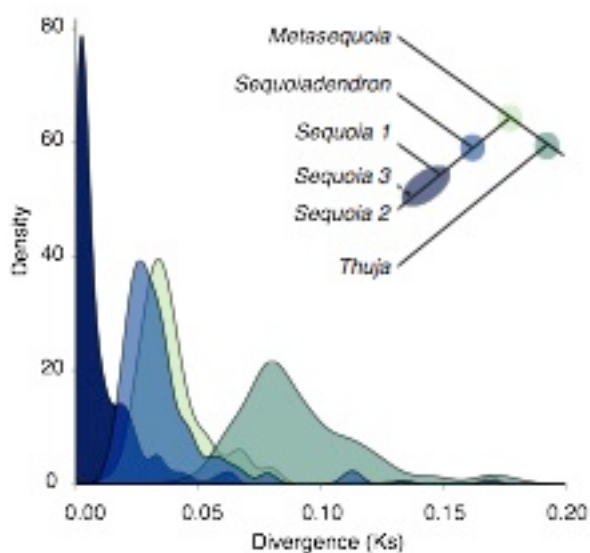
242     Californian redwood clade, we estimated phylogenetic trees for all genes with more than

243     one sequence variant in *Sequoia*.  A total of 217 genes were present in two or three copies

244     in *S. sempervirens*. The optimal tree for 186 of these alignments (85.7%) showed

245     monophyly of the *S. sempervirens* copies with *Sequoia* sister to *Sequoiadendron* (Fig. 3),

246     with 97% of these trees well-supported (i.e., having a bootstrap > 0.70). The remaining

247     31 genes (14%) either contradicted monophyly of *S. sempervirens* copies, supporting

248     several other possible relationships, or lacked clear resolution of species relationships.

249

250     Based on ML estimates using a codon model in PAML, we could calculate the patristic

251     Ka and Ks distances between each pair of tips for on each genes tree. Doing this on the

252     176 well-supported gene trees that yielded a monophyletic *Sequoia*, average phylogenetic

253     $K_s$ among *Sequoia* gene copies was 0.013. This was approximately one-third of the Ks

254     separating *Sequoia* sequences from other redwoods (Figure 4).

255

256     **Figure 4: Tree-based divergence estimates in Ks**
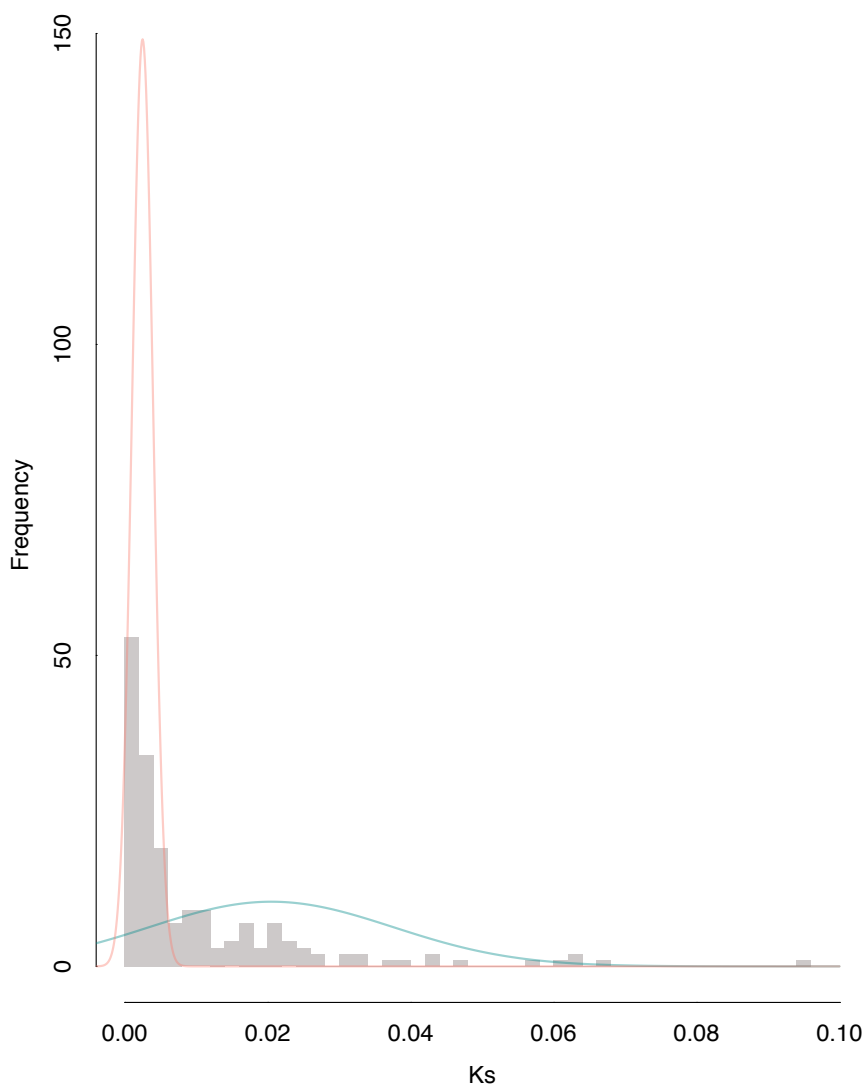


257

258     Density distribution of divergence estimates (in Ks). For Distributions are colored to indicate

259     corresponding nodes on the tree.

260

261

262

263

264    Figure 5: Age distribution of *Sequoia* variants. Colored lines denote normal distributions fit with EMMIX.

265

266    We tested whether the patristic Ks estimates between *S. sempervirens* copies are sampled

267    from one or two normal distributions. If hexaploidy arose from two sequential WGD

268    events, there should be two, distinct normal distributions. We used EMMIX to fit a

269    mixture model of normal distributions to the PAML Ks estimates. Based on AIC and BIC

270    scores, the presence of two Gaussian distributions provides a better fit to the $K_s$ distance

271    data. Figure 5 shows the best fitting pair of distributions. Although it is difficult to

272    reliably translate $K_s$ into absolute age, using a generic average mutation rate for conifers

273    of $0.68 \times 10^{-9}$ synonymous substitutions per site per year (Buschiazzo et al., 2012), these

274    peaks correspond to ~3 Ma and 10 Ma.

275

276    **DISCUSSION**

277

278    **Transcriptome sequencing in the redwoods supports a sister group relationship**

279    **between *Sequoia* and *Sequoiadendron*.**

280    Bayesian concordance analysis of single copy genes overwhelmingly supports

281    *Sequoiadendron* as the closest relative of *Sequoia*. This conclusion is in agreement with

282    decades of previous work based on morphology, karyotype, and chloroplast sequence

283    data (e.g. Brunsfield et al., 1994; Gadek et al., 2000; Kusumi et al., 2001).

284

285    We found genes supporting two minor topologies, one with a *Sequoia-Metasequoia* clade

286    and the other with a *Sequoiadendron–Metasequoia* clade. These discordant topologies

287    could be due to incomplete lineage sorting (ILS), which arises when multiple gene copies

288    (or alleles) persist between sequential splits in a population tree. In this case, the two

289    minor trees have similar concordance factors, 0.010 and 0.11, and their associated

290    credibility intervals overlap.  This pattern is consistent with ILS, which predicts that the

291    alternative minor topologies should have equal CFs (Baum 2007). Furthermore, given a

292    concordance factor of 0.80, coalescent theory would predict that *Sequoia*-

293    *Sequoiadendron* clade is subtended by a population lineage whose duration was ~1.22 Ne

294    generations, where Ne is the effective population size (Allman et al., 2011; Larget et al.

295    2011). However, it is also possible that the internal branch is considerably longer and

296    discordance is due to other factors such as mistaken orthology. The fact that the two

297    minor histories have similar concordance factors tends to argue against introgression or

298    hybridization as an important phenomenon in the group.

299

300    **Hexaploidy in *Sequoia* did not involve hybridization among extant redwood**

301    **lineages.**

302    Our phylogenetic results support an autopolyploid origin for hexaploid *Sequoia*, with no

303    evidence to support hybridization among modern redwood lineages. Single-copy trees

304    convey strong support for *Sequoiadendron* as the closest relative of *Sequoia*, suggesting

305    there was no genome contribution from *Metasequoia*. The lack of evidence that

306    *Metasequoia* was involved with the polyploid origins of *Sequoia* puts some long-held

307    hypotheses to rest (e.g. Stebbins, 1948; Saylor & Simons, 1970). However, as these

308    phylogenies include only one copy for hexaploid *Sequoia*, they could not distinguish

309    between autopolyploidy within the *Sequoia* lineage or autoallopolyploidy within the

310    *Sequoiadendron*-*Sequoia* clade. Single-copy trees may also be inconclusive due to

311    extreme copy-specific expression or genome dominance, where genes from one parental

312    genome are preferentially expressed (e.g. Woodhouse et al., 2014). Therefore, we sought

313    additional evidence by studying orthogroups that included 2 or 3 distinct sequence

314    variants, putatively homeologs, from *Sequoia*. Phylogenetic analyses of these

315    orthogroups strongly support monophyly of *Sequoia* homeologs, suggesting that all gene

316    copies in *Sequoia* originate from the same redwood lineage.

317

318    **Polyploidy in *Sequoia* arose relatively recently**.

319    The similarity of the $K_s$ plots obtained from polyploid *Sequoia* and diploids

320    *Sequoiadendron* and *Metasequoia* (Fig. 2), and specifically the lack of a recent peak

321    restricted to *Sequoia*, is initially surprising, as these methods have been widely used to

322    diagnose polyploidization events in numerous plant lineages (e.g. Barker et al., 2008; Jiao

323    et al., 2011). This pattern might be expected if autopolyploidy had occurred very

324    recently, such that the level of divergence among homeologs is not much different than

325    that among alleles at a particular locus (Vanneste et al. 2013), but the fossil data suggests

326    polyploidization as early as the Eocene. One possible explanation for the lack of a

327    polyploidization peak is that only one homeolog is expressed in leaves. Such genome

328    dominance has been observed in other polyploid species (e.g., Adams et al. 2004).

329    However, the fact that we found many genes with two or three distinct copies in *Sequoia*

330    but only one in each diploid argues against uniform silencing of all but one homeolog.

331

332    To further explore the history of gene duplication, we inferred trees for alignments that

333    included one transcript in diploids and two or three from *Sequoia* and then inferred the

334    branch lengths of this tree in $K_s$ units. We found that $K_s$ estimates between even the most

335    divergent *Sequoia* homeologs were very low (>0.10). One possible explanation is that

336    *Sequoia* experienced a long period of multisomic inheritance following autopolyploidy

337  during which time homeologs tended to be repeatedly recombined, resulting in much

338  lower $K_s$ values (described in Wolfe, 2001). These observations highlight some caveats

339  of using paralog age distribution graphs alone to infer recent polyploidization events, or

340  to study ancient whole genome duplication events that were accompanied by extended

341  periods of multisomic inheritance.

342

343  Fitting a mixture model of normal distributions to $K_s$ estimates between homeologs

344  yielded two distinct, but overlapping Gaussian distributions. This suggests two whole

345  genome duplication events are included in our age distribution data.  Using a mutation

346  rate calibration for conifer Ks divergence, we estimated the timing of the first whole

347  genome duplication in *Sequoia* to have occurred around 10 Ma, with the second

348  occurring more recently, about 3 Ma. These dates are in apparent contradiction to the

349  discovery of *Sequoia* fossils in the Eocene (33-53 Ma) with guard cells of a size taken to

350  be indicative of polyploidy (Ma et al., 2005). One possible explanation for this

351  discrepancy is that the mutation rate is three-fold lower in *Sequoia*  (or redwoods in

352  general) than in other conifers. However, although some redwoods may have extremely

353  long life spans, such a great different in the rate of synonymous substitutions seems

354  improbable.

355

356  A second possibility is that the Eocene fossils represent an independent instance of

357  polyploidy in a closely related lineage that was misclassified as being in *Sequoia*. It is

358  noteworthy that some plant groups that acquire the propensity to undergo polyploidy, do

359  so repeatedly, a possible case in point being the *Ephedra* lineage*,* which appears to have

360  experienced multiple whole genome duplication events (Ickert-Bond, 2003). Further

361  evaluating this hypothesis would require measurements of guard cells in a much larger

362  number of different aged *Sequoia* fossils from different geographic locations.

363

364  The final possible explanation for the low divergence of putative homeologs in *Sequoia* is

365  that while autopolyploidy occurred in the Eocene (or even earlier), multisomic

366  inheritance persisted for a long period of time, possibly even to the present for some loci.

367  In such a case the gene duplication events we dated would not correspond to the

368 polyploidy event per se but would reflect subsequent, recombinational homogenization.

369 This hypothesis is consistent with multivalent formation in modern *Sequoia*, and suggests

370 a very slow diploidization process following whole genome duplication in *Sequoia*.

371

372 **Implications for polyploidization patterns in gymnosperms**.

373 Given what we know about polyploidy in *Sequoia*, what conclusions can we draw about

374 patterns of polyploidization in gymnosperms overall? With the exception of *Ephedra*,

375 instances of polyploid gymnosperms are limited to monospecific genera (e.g. *Sequoia*,

376 *Fitzroya*), or even just to polyploid individuals within diploid species (e.g. *Juniperus x*

377 *pftizeriana*; Ahuja, 2005). If polyploidy in gymnosperms is associated with small clades,

378 as seems to be the case, we can infer that polyploidy either hinders speciation or

379 promotes extinction of gymnosperm lineages, or both.

380

381 The apparent mismatch between the inferred age of gene duplication and the timing of

382 polyploidization as seen in the fossil record suggests an intriguing hypothesis to explain

383 the paucity of polyploidy in gymnosperms. Perhaps diploidization happens more slowly

384 in gymnosperms (except perhaps *Ephedra*) than in angiosperms. The main long-term

385 benefits of polyploidy (potential sub- and neo-functionalization of genes) require

386 divergence among homeologous chromosomes, which can only happen once loci are

387 diploidized. Thus, continued multisomic inheritance precludes the emergence of any

388 evolutionary advantage in polyploid lineages.

389

390 If polyploidy in gymnosperms is more burden than boon, the persistence of hexaploid

391 *Sequoia* may reflect an ability to avoid extinction rather than superior fitness. In this

392 regard it is perhaps noteworthy that *S. sempervirens* manifests some traits that might help

393 stave of extinction, namely clonal reproduction, self-compatibility, and extreme

394 longevity. In coast redwood populations, suckers often emerge from the base of adult

395 trees, extending generation time (meiosis-to-meiosis) almost indefinitely. Furthermore,

396 production of asexual stands may lead to abundant genetic selfing among clonal ramets,

397 as coast redwoods are self-compatible (Burns & Honkala, 1990). This means that a

398 spontaneous polyploid, perhaps gaining the transient advantage of fixed heterozygosity,

399    could spread by a combination of asexual reproduction and selfing. It is conceivable,

400    therefore, that even after the erosion of fixed heterozygosity the lineage could persist

401    despite never gaining the long-term advantages typically associated with polyploidy,

402    instead suffering the concomitant problem of enlarged genome size. The only other

403    natural polyploid in Cupressaceae, *Fitzroya cupressoides*, is a putative autotetraploid.

404    Like *Sequoia*, *Fitzroya* is both long-lived and capable of clonal reproduction (Silla et al.,

405    2002).  Thus, while more work is needed to evaluate the occurrence of multisomic

406    inheritance in both polyploid species (e.g. *Sequoia, Fitzroya)* and polyploid clones

407    *Juniperus x pftizeriana*, our hypothesis can both explain the rarity of neopolyploidy in

408    gymnosperms and why *Sequoia* is an exception to this general rule.

409

418    **AUTHOR CONTRIBUTIONS**

419    ADS and DB designed the research and wrote the manuscript, ADS collected the data,

420    ADS and NS analyzed the data.

421

422    **REFERENCES**

423

424    Adams, K. L., Percifield, R., & Wendel, J. F. (2004). Organ-specific silencing of

425    duplicated genes in a newly synthesized cotton allotetraploid. *Genetics*, *168*(4), 2217-

426    2226.

427

428    Adams, K. L., & Wendel, J. F. (2005). Polyploidy and genome evolution in

429    plants. *Current opinion in plant biology*, *8*(2), 135-141.

430

431   Ahuja, M. R. (2005). Polyploidy in gymnosperms: revisited. *Silvae Genetica,54*(2), 59-

432   68.

433

434   Ahuja, M. R. (2009). Genetic constitution and diversity in four narrow endemic redwoods

435   from the family Cupressaceae. *Euphytica, 165*(1), 5-19.

436

437   Ahuja, M. R., & Neale, D. B. (2002). Origins of polyploidy in coast redwood (Sequoia

438   sempervirens (D. Don) Endl.) and relationship of coast redwood to other genera of

439   Taxodiaceae. *Silvae Genetica, 51*(2-3), 93-99.

440

441   Ahuja, M. R., & Neale, D. B. (2005). Evolution of genome size in conifers. *Silvae*

442   *genetica, 54*(3), 126-137.

443

444   Allman, E. S., Degnan, J. H., & Rhodes, J. A. (2011). Identifying the rooted species tree

445   from the distribution of unrooted gene trees under the coalescent.*Journal of mathematical*

446   *biology, 62*(6), 833-862.

447

448   Ané, C., Larget, B., Baum, D. A., Smith, S. D., & Rokas, A. (2007). Bayesian estimation

449   of concordance among gene trees. *Molecular Biology and Evolution,24*(2), 412-426.

450

451   Barker, M. S., Kane, N. C., Matvienko, M., Kozik, A., Michelmore, R. W., Knapp, S. J.,

452   & Rieseberg, L. H. (2008). Multiple paleopolyploidizations during the evolution of the

453   Compositae reveal parallel patterns of duplicate gene retention after millions of

454   years. *Molecular Biology and Evolution, 25*(11), 2445-2455.

455

456   Baum, D. A. (2007). Concordance trees, concordance factors, and the exploration of

457   reticulate genealogy. *Taxon*, 417-426.

458

459   Beaulieu, J. M., Leitch, I. J., Patel, S., Pendharkar, A., & Knight, C. A. (2008). Genome

460   size is a strong predictor of cell size and stomatal density in angiosperms. *New*

461   *Phytologist*, *179*(4), 975-986.

462

463   Brunsfeld, S. J., Soltis, P. S., Soltis, D. E., Gadek, P. A., Quinn, C. J., Strenge, D. D., &

464   Ranker, T. A. (1994). Phylogenetic relationships among the genera of Taxodiaceae and

465   Cupressaceae: evidence from rbcL sequences. *Systematic Botany*, 253-262.

466

467   Burns RM, Honkala BH: Silvics of North America. In *Agriculture Handbook 654*. 2nd

468   edition. Washington, DC: U.S: Department of Agriculture, Forest Service; 1990

469

470   Chaudhuri, P., & Marron, J. S. (1999). SiZer for exploration of structures in

471   curves. *Journal of the American Statistical Association*, *94*(447), 807-823.

472

473   Douhovnikoff, V., Cheng, A. M., & Dodd, R. S. (2004). Incidence, size and spatial

474   structure of clones in second-growth stands of coast redwood, Sequoia sempervirens

475   (Cupressaceae). *American Journal of Botany*, *91*(7), 1140-1146.

476

477   Douhovnikoff, V., & Dodd, R. S. (2011). Lineage divergence in coast redwood (Sequoia

478   sempervirens), detected by a new set of nuclear microsatellite loci.*The American Midland*

479   *Naturalist*, *165*(1), 22-37.

480

481   Doyle, J. (1945). Naming of the redwoods. *Nature*, *155*, 254-257.

482

483   Eckenwalder, J. E. (1976). Re-evaluation of Cupressaceae and Taxodiaceae: a proposed

484   merger. *Madrono*, *23*(5), 237-256.

485

486   Edgar, R.C. (2004) MUSCLE: multiple sequence alignment with high accuracy and high

487   throughput. *Nucleic Acids Res.* **32**(5):1792-1797

488

489    Edgar, R.C. (2004) MUSCLE: a multiple sequence alignment method with reduced time

490    and space complexity. *BMC Bioinformatics*, (**5**) 113

491

492    Fozuar, B. S., & Libby, W. J. (1968). Chromosomes of Sequoia sempervirens; 8-

493    hydroxy-quinoline-castor oil pretreatment for improving preparation.*Biotechnic &*

494    *Histochemistry*, *43*(2), 97-100.

495

496    Gadek, P. A., Alpers, D. L., Heslewood, M. M., & Quinn, C. J. (2000). Relationships

497    within Cupressaceae sensu lato: a combined morphological and molecular

498    approach. *American Journal of Botany*, *87*(7), 1044-1057.

499

500    Grabherr MG, Haas BJ, Yassour M, Levin JZ, Thompson DA, Amit I, Adiconis X, Fan

501    L, Raychowdhury R, Zeng Q, Chen Z, Mauceli E, Hacohen N, Gnirke A, Rhind N, di

502    Palma F, Birren BW, Nusbaum C, Lindblad-Toh K, Friedman N, Regev A. Full-length

503    transcriptome assembly from RNA-seq data without a reference genome.

504

505    Haas, B. J., Papanicolaou, A., Yassour, M., Grabherr, M., Blood, P. D., Bowden, J., ... &

506    Regev, A. (2013). De novo transcript sequence reconstruction from RNA-seq using the

507    Trinity platform for reference generation and analysis.*Nature protocols*, *8*(8), 1494-1512.

508

509    Huelsenbeck, J. P. and F. Ronquist. 2001. MRBAYES: Bayesian inference of phylogeny.

510    Bioinformatics 17:754-755.

511

512    Ickert-Bond, S. M. 2003. Systematics of New World *Ephedra*L. (Ephedraceae):

513    integrating morphological and molecular data. Ph.D. dissertation, Arizona State

514    University, Tempe, Arizona, USA

515

516    Ishii, H. R., Azuma, W., Kuroda, K., & Sillett, S. C. 2014. Pushing the limits to tree

517    height: could foliar water storage compensate for hydraulic constraints in Sequoia

518    sempervirens?. Functional Ecology, 28(5), 1087-1093.

519

520 Jiao, Y., Wickett, N. J., Ayyampalayam, S., Chanderbali, A. S., Landherr, L., Ralph, P.

521 E., ... & Leebens-Mack, J. (2011). Ancestral polyploidy in seed plants and

522 angiosperms. *Nature*, *473*(7345), 97-100.

523

524 Kent, W. J. (2002). BLAT—the BLAST-like alignment tool. *Genome research*,*12*(4),

525 656-664.

526

527 Khoshoo, T. N. (1959). Polyploidy in gymnosperms. *Evolution*, 24-39.

528

529 Kusumi, J., Tsumura, Y., Yoshimaru, H., & Tachida, H. (2000). Phylogenetic

530 relationships in Taxodiaceae and Cupressaceae sensu stricto based on matK gene, chlL

531 gene, trnL-trnF IGS region, and trnL intron sequences. *American Journal of*

532 *Botany*, *87*(10), 1480-1488.

533

534 Larget, B. R., Kotha, S. K., Dewey, C. N., & Ané, C. (2010). BUCKy: gene tree/species

535 tree reconciliation with Bayesian concordance analysis.*Bioinformatics*, *26*(22), 2910-

536 2911.

537

538 Lechner, M., Findeiß, S., Steiner, L., Marz, M., Stadler, P. F., & Prohaska, S. J. (2011).

539 Proteinortho: Detection of (Co-) orthologs in large-scale analysis. *BMC*

540 *bioinformatics*, *12*(1), 124.

541

542 Leslie, A. B., Beaulieu, J. M., Rai, H. S., Crane, P. R., Donoghue, M. J., & Mathews, S.

543 (2012). Hemisphere-scale differences in conifer evolutionary dynamics. *Proceedings of*

544 *the National Academy of Sciences*, *109*(40), 16217-16221.

545

546 Ma, Q. W., Li, F. L., & Li, C. S. (2005). The coast redwoods (Sequoia, Taxodiaceae)

547 from the Eocene of Heilongjiang and the Miocene of Yunnan, China. *Review of*

548 *Palaeobotany and Palynology*, *135*(3), 117-129.

549

550    Mao, K., Milne, R. I., Zhang, L., Peng, Y., Liu, J., Thomas, P., ... & Renner, S. S. (2012).

551    Distribution of living Cupressaceae reflects the breakup of Pangea.*Proceedings of the*

552    *National Academy of Sciences*, *109*(20), 7793-7798.

553

554    Mclachlan G, Peel D, Basford K, Adams P. 1999. The EMMIX software for the fitting of

555    mixtures of normal and t- components. J Stat Softw. 4:2.

556    Miki, S., & Hikita, S. (1951). Probable chromosome number of fossil Sequoia and

557    Metasequoia found in Japan. *Science*, *113*(2923), 3-4.

558

559    Otto, S. P., & Whitton, J. (2000). Polyploid incidence and evolution. *Annual review of*

560    *genetics*, *34*(1), 401-437.

561

562    Parisod, C., Holderegger, R., & Brochmann, C. (2010). Evolutionary consequences of

563    autopolyploidy. *New Phytologist*, *186*(1), 5-17.

564

565    R Core Team (2013). R: A language and environment for statistical   computing. R

566    Foundation for Statistical Computing, Vienna, Austria.   URL http://www.R-project.org/.

567

568    Ramsey, J., & Schemske, D. W. (2002). Neopolyploidy in flowering plants.*Annual*

569    *review of ecology and systematics*, 589-639.

570

571    Rogers, D. L. (1997). Inheritance of allozymes from seed tissues of the hexaploid

572    gymnosperm, Sequoia sempervirens(D. Don) Endl.(Coast redwood). *Heredity*, *78*(2),

573    166-175.

574

575    Rogers, D. L. (2000). Genotypic diversity and clone size in old-growth populations of

576    coast redwood (Sequoia sempervirens). *Canadian Journal of Botany*, *78*(11), 1408-1419.

577

578    Ronquist, F. and J. P. Huelsenbeck. 2003. MRBAYES 3: Bayesian phylogenetic

579    inference under mixed models. Bioinformatics 19:1572-1574.

580

581   Saylor, L. C., & Simons, H. A. (1970). Karyology of Sequoia sempervirens: karyotype
582   and accessory chromosomes. *Cytologia*, *35*(2), 294-303.
583
584   Schlarbaum, S. E., & Tsuchiya, T. (1984). Cytotaxonomy and phylogeny in certain
585   species of Taxodiaceae. *Plant systematics and evolution*, *147*(1-2), 29-54.
586
587   Schlarbaum, S. E., Tsuchiya, T., & Johnson, L. C. (1984). The chromosomes and
588   relationships of Metasequoia and Sequoia (Taxodiaceae): an update.*Journal of the Arnold*
589   *Arboretum*, *65*(2), 251-254.
590
591   Schulz, C., & Stützel, T. (2007). Evolution of taxodiaceous Cupressaceae
592   (Coniferopsida). *Organisms Diversity & Evolution*, *7*(2), 124-135.
593
594   Silla, F., Fraver, S., Lara, A., Allnutt, T. R., & Newton, A. (2002). Regeneration and
595   stand dynamics of Fitzroya cupressoides (Cupressaceae) forests of southern Chile's
596   Central Depression. *Forest Ecology and Management*, *165*(1), 213-224.
597
598   Stamatakis, A. (2006). RAxML-VI-HPC: maximum likelihood-based phylogenetic
599   analyses with thousands of taxa and mixed models.*Bioinformatics*, *22*(21), 2688-2690.
600
601   Stebbins, G. L. (1947). Types of polyploids: their classification and
602   significance.*Advances in genetics*, *1*, 403-29.
603
604   Stebbins, G. L. (1948). The chromosomes and relationships of Metasequoia and
605   Sequoia. *Science*, *108*(2796), 95-98.
606
607   Vanneste, K., Van de Peer, Y., & Maere, S. (2013). Inference of genome duplications
608   from age distributions revisited. *Molecular biology and evolution,30*(1), 177-190.
609
610   Wolfe, K. H. (2001). Yesterday's polyploids and the mystery of diploidization.*Nature*
611   *Reviews Genetics*, *2*(5), 333-341.

612

613   Woodhouse, M. R., Cheng, F., Pires, J. C., Lisch, D., Freeling, M., & Wang, X. (2014).

614   Origin, inheritance, and gene regulatory consequences of genome dominance in

615   polyploids. *Proceedings of the National Academy of Sciences*,*111*(14), 5283-5288.

616

617   Yang, Z. 1997. PAML: a program package for phylogenetic analysis by maximum

618   likelihood

619   Computer Applications in BioSciences 13:555-556.

620

621   Yang, Z. 2007. PAML 4: a program package for phylogenetic analysis by maximum

622   likelihood. Molecular Biology and Evolution 24: 1586-1591

623

624   Yang, Z. Y., Ran, J. H., & Wang, X. Q. (2012). Three genome-based phylogeny of

625   Cupressaceae sl: Further evidence for the evolution of gymnosperms and Southern

626   Hemisphere biogeography. *Molecular phylogenetics and evolution*,*64*(3), 452-470.

627

628   Zhang Z, Li J, Zhao XQ, Wang J, Wong GK, Yu J: KaK$_s$ Calculator: Calculating Ka and

629   K$_s$ through model selection and model averaging.*Genomics Proteomics*

630   *Bioinformatics* 2006 , 4:259-263.

631

632