

1 **Convergent evolution of polyploid genomes from across the eukaryotic tree of life**

2
3 Yue Hao^{1,*}, Jonathon Fleming^{2,*}, Joanna Petterson³, Eric Lyons⁴, Patrick P. Edger^{5,6}, J. Chris
4 Pires⁷⁻⁹, Jeffrey L. Thorne^{2,10-12}, and Gavin C. Conant^{2,10,12,†}

5
6 ¹Biodesign Center for Mechanisms of Evolution, Arizona State University, Tempe, AZ, U.S.A.

7 ²Bioinformatics Research Center and ³Department of Biomedical Engineering, North Carolina
8 State University, Raleigh, NC, U.S.A.

9 ⁴School of Plant Sciences, University of Arizona, Tucson AZ, U.S.A.

10 ⁵Department of Horticulture, Michigan State University, East Lansing, MI, U.S.A.

11 ⁶Ecology, Evolutionary Biology and Behavior, Michigan State University, East Lansing, MI,
12 U.S.A.

13 ⁷International Plant Science Center; New York Botanical Garden, Bronx, NY, U.S.A.

14 ⁸Division of Biological Sciences, University of Missouri-Columbia, MO, U.S.A.

15 ⁹Bond Life Sciences Center, University of Missouri-Columbia, MO, U.S.A.

16 ¹⁰Program in Genetics, ¹¹Department of Statistics and ¹²Department of Biological Sciences,
17 North Carolina State University, Raleigh, NC, U.S.A.

18 *These authors contributed equally to this work

19
20 †Correspondence: G. Conant, Campus Box 7566, NCSU Campus, Raleigh, NC, 27695, U. S.
21 A., gconant@ncsu.edu

22
23
24 Running Head: Convergent patterns of evolution after polyploidy

25
26 Keywords: polyploidy, convergent evolution, reciprocal gene loss, evolutionary model

27 **Abstract:**

28 By modeling the homoeologous gene losses that occurred in fifty genomes deriving from ten
29 distinct polyploidy events, we show that the evolutionary forces acting on polyploids are
30 remarkably similar, regardless of whether they occur in flowering plants, ciliates, fishes or
31 yeasts. We show that many of the events show a relative rate of duplicate gene loss prior to the
32 first post-polyploidy speciation that is significantly higher than in later phases of their evolution.
33 The relatively weak selective constraint experienced by the single-copy genes these losses
34 produced leads us to suggest that most of the purely selectively neutral duplicate gene losses
35 occur in the immediate post-polyploid period. Nearly all of the events show strong evidence of
36 biases in the duplicate losses, consistent with them being allopolyploidies, with two distinct
37 progenitors contributing to the modern species. We also find ongoing and extensive reciprocal
38 gene losses (RGL; alternative losses of duplicated ancestral genes) between these genomes. With
39 the exception of a handful of closely related taxa, all of these polyploid organisms are separated
40 from each other by tens to thousands of reciprocal gene losses. As a result, it is very unlikely that
41 viable diploid hybrid species could form between these taxa, since matings between such hybrids
42 would tend to produce offspring lacking essential genes. It is therefore possible that the relatively
43 high frequency of recurrent polyploidies in some lineages may be due to the ability of new
44 polyploidies to bypass RGL barriers.

45

46 **Introduction**

47 That organisms with doubled genomes existed was evident early in the history of genetics
48 (KUWADA 1911; CLAUSEN AND GOODSPEED 1925), and a lively debate was entered as to the
49 implications of this fact. Wagner (1970) declared polyploidy to be “evolutionary noise” the same
50 year that Susumu Ohno (1970) was giving it pride of place among the forces generating
51 evolutionary innovations. The advent of genome sequencing changed the ground of this debate,
52 opening new horizons of time for studies of the prevalence and influence of polyploidy. We
53 know now that great branches of the eukaryotic evolutionary tree, including the vertebrates, all
54 flowering plants and many yeasts, descend from ancient polyploids (VAN DE PEER *et al.* 2017),
55 events that were difficult or impossible to detect with older data. For reasons that are not yet
56 fully understood, many of these groups also show recurrent polyploidies, especially flowering
57 plants (SOLTIS *et al.* 2009) and teleost fishes (BRAASCH AND POSTLETHWAIT 2012).

58 With this extensive new set of polyploidies as a resource, other old questions can also be
59 revisited, such as the relative prevalence of auto- and allopolyploids (STEBBINS JR 1947).
60 Allopolyploidy refers to hybridizations between distinct species that result in doubled (or more)
61 genomes, while autopolyploids are derived from a single progenitor species (KUWADA 1911;
62 CLAUSEN AND GOODSPEED 1925; STEBBINS JR 1947). Analyses of several paleopolyploid
63 genomes have shown that while gene losses are common after polyploidy, in many cases the
64 losses are not experienced equally by the two parental subgenomes (THOMAS *et al.* 2006; EMERY
65 *et al.* 2018), a pattern known as biased fractionation. These biases are plausible but not definitive
66 indicators of allopolyploidy.

67 There has also been controversy as to whether and how polyploidy affects the rate of
68 speciation. Werth and Windham (1991) proposed that reciprocal loss of expression at duplicated
69 loci could create Bateson–Dobzhansky–Muller incompatibilities between populations (see ORR
70 1996 for a history of this concept), because matings between them would give rise to offspring
71 that did not express either copy of the gene. Reciprocal gene losses (RGLs) after polyploidy are
72 an example of this process, and, were those genes essential, the offspring lacking their presence
73 or expression would be inviable (WERTH AND WINDHAM 1991) Such incompatibilities have been
74 observed both in the wild and the laboratory (MIZUTA *et al.* 2010; MACLEAN AND GREIG 2011).
75 Muir and Hahn (2015) emphasize that RGL requires a period of reproductive isolation to form.

76 In the case of the ancient polyploidy in bakers' yeast and its relatives, RGLs are
77 commonly found between the descendant genomes, suggesting the potential for polyploidy to
78 create new species by purely neutral means (SCANNELL *et al.* 2006; SCANNELL *et al.* 2007).
79 However, direct analyses of the speciation and extinction rates of lineages with and without
80 recent polyploidy events has yielded inconclusive results, with some studies claiming reduced
81 net diversification rates among polyploids and others disagreeing (MAYROSE *et al.* 2011; SOLTIS
82 *et al.* 2014a). More generally, the immediate and long-term adaptive value of polyploidy remains
83 unclear: for instance, allopolyploids combine hybridizations with genome doubling and may
84 derive immediate advantages from the hybridization effects rather than the doubling itself
85 (SOLTIS *et al.* 2014b). Increased stress tolerance in polyploid organisms due to a variety of
86 immediate and evolutionary mechanisms (VAN DE PEER *et al.* 2021) has also been invoked to
87 argue for a radiation of polyploidy coincident with global catastrophes such as the KT mass
88 extinction (FAWCETT *et al.* 2009).

89 Hence, while many studies of the resolution of individual polyploidies have been made
90 (MAERE *et al.* 2005; SCANNELL *et al.* 2006; THOMAS *et al.* 2006; BUGGS *et al.* 2009a;
91 WOODHOUSE *et al.* 2010; BRAASCH AND POSTLETHWAIT 2012) and a few comparisons of several
92 events are available (PATERSON *et al.* 2006; DE SMET *et al.* 2013; GARSMEUR *et al.* 2013; EMERY
93 *et al.* 2018), no deep, cross-kingdom analyses of the patterns of post-polyploidy evolution using
94 uniform and rigorous models have been undertaken. In the same vein, the similarities in which
95 types of homoeologous genes are retained and lost after polyploidy (SEOIGHE AND WOLFE 1999;
96 BLANC AND WOLFE 2004; PATERSON *et al.* 2006; FREELING 2009; DE SMET *et al.* 2013), as well
97 as the prevalence of biased fractionation (THOMAS *et al.* 2006; GARSMEUR *et al.* 2013; EMERY *et al.*
98 *et al.* 2018) are examples of pattern-based convergent evolution (STAYTON 2015). However, a
99 broad phylogenomic analysis of polyploidy with such uniform models is needed to ground this
100 qualitative description of convergence with estimates of how similar or different the model
101 parameters describing duplicate retention or biased fractionation are across these polyploidy
102 events.

103 Using our tool for modeling the evolution of polyploid genomes, POInT (the Polyploidy
104 Orthology Inference Tool; CONANT AND WOLFE 2008), we explored the resolution of ten
105 independent polyploidies. We adopt the term “homoeolog” below to refer to homologous genes
106 produced by any type of polyploidy rather than “duplicate” or “ohnolog” because the events

107 considered comprise several distinct types of polyploidy. The hallmark of polyploidy in a
108 genome is a pattern of interleaved synteny, comprising not just the surviving homoeologs but
109 also single-copy genes that are now found in interleaved positions on pairs (or more) of
110 chromosomal segments homologous to the ancestral single-copy regions. In Figure 1A, we show
111 an example of this evolutionary process, which yields conserved synteny blocks in the extant
112 genomes. Those synteny blocks differ between genomes, meaning it is necessary to “phase”
113 them into orthologous regions. As shown in Figure 1B, for a set of n tetraploid genomes, there
114 are 2^n possible orthology relationships at each ancestral locus. We use the term “pillar” to denote
115 all of the genes or lost homoeologs at such a locus. POInT computes the likelihood of the
116 observed homoeolog presence/absence data at each pillar for each possible orthology
117 relationship. Via a hidden Markov model (HMM) that combines the possible orthology
118 relationships for each pillar with the syntenic organization among pillars (Figure 1C), POInT
119 employs posterior decoding to infer orthology estimates for each pillar with associated posterior
120 probabilities (top of Figure 1D) as well as estimates of the model parameters describing the
121 process of homoeolog loss (Figure 2).

122 Our analyses here encompass a total of 50 polyploid genomes and more than 460,000
123 individual genes (Figure 3). We find that the patterns of gene loss after these different events
124 show strikingly similar patterns, with strong evidence for biased fractionation and homoeolog
125 fixation. Using synonymous substitutions as an evolutionary clock, we show that the rate of gene
126 loss immediately after polyploidy is generally higher than in later periods. RGL is also prevalent
127 after all of these polyploidy events, and we suggest it might introduce barriers to hybridization
128 that could be overcome through subsequent allopolyploidy events.

129

130 **Methods**

131 *Synteny block inference.*

132 Our three-step pipeline for inferring blocks of pillars with n -fold conserved synteny
133 (NCS) produced by polyploidy (CONANT 2020) first uses GenomeHistory (CONANT AND
134 WAGNER 2002) to find all pairs of homologous genes between each polyploid genome and a
135 outgroup lacking the event in question (see Supplemental Table 1 for genome details and
136 Supplemental Table 2 for parameters). The second step seeks to place these homologous genes
137 into $N:1$ relationships between the polyploid genome and the outgroup ($N=2$ for a WGD, $N=3$

138 for a hexaploidy and $N=4$ for an octoploidy). Using simulated annealing (KIRKPATRICK *et al.*
 139 1983), this step proposes sets of ordered pillars, each of which contains a single gene from the
 140 outgroup that lacks the polyploidy (G) and no more than N of the homologs of that gene from the
 141 polyploid genome. The annealing algorithm then seeks a combination of these assignments and a
 142 relative ordering of the m outgroup genes $G_1..G_m$ that maximizes the number of synteny
 143 relations. We define two genes to be in synteny if they are neighbors in the genome, ignoring any
 144 genes without homologs to the compared genome. In the third step, these NCS blocks for each
 145 polyploid genome are merged across all of the polyploid genomes. In this merging, only pillars
 146 where we have at least one homologous and syntenic gene from each polyploid genome are
 147 included. With the set of merged pillars, a further simulated annealing search is undertaken to
 148 infer a global pillar order that minimizes the number of synteny breaks. While not strictly an
 149 ancestral genome inference (SANKOFF AND BLANCHETTE 1998), it is helpful to think of this
 150 optimal ordering as approximating the order of the genes just prior to the polyploidy event. Our
 151 previous work has shown that this inference approach is highly specific, with no apparent cases
 152 of paralogous genes not created by the polyploidies in question being included in the pillars
 153 (EMERY *et al.* 2018; CONANT 2020).

154

155 *Modeling polyploidies with POInT*

156 At each pillar, POInT calculates the probability of the observed gene presence-absence
 157 data conditional upon all possible orthology relationships and a phylogeny. It carries this
 158 uncertainty in orthology through its likelihood computations using a hidden Markov model that
 159 resembles the Lander-Green approach for constructing linkage maps on a pedigree (LANDER AND
 160 GREEN 1987). The parameter θ_i corresponds to the probability that the inferred orthology
 161 relationships change between syntenic neighbors at pillars $i-1$ and i . When a pair of pillars are
 162 separated by a synteny break (i.e., the two genes are not each other's chromosomal neighbors),
 163 their orthology relationships are independent (i.e., $\theta_i=1/2$). Otherwise, $\theta_i=\theta$, a global parameter
 164 estimated from the data by maximum likelihood.

165 This modeling framework allows for testing hypotheses about post-polyploidy gene
 166 losses. We have extensively validated it in several prior contributions (CONANT AND WOLFE
 167 2008; CONANT 2014; EMERY *et al.* 2018; CONANT 2020). For tetraploidies, we analyzed three
 168 phenomena: fixation of homoeolog pairs, biased fractionation and overly frequent parallel losses

169 of the same homoeolog on independent branches of the phylogeny (Supplemental Table 3). For
170 the *Brassica* hexaploidy and nematode triplication events, we focused on differences in
171 homoeolog loss rates between the three subgenomes (Supplemental Table 4). We further allowed
172 the root branch to have separate values of the model parameters to account for the two-step
173 nature of hexaploidy formation (Figure 2; TANG *et al.* 2012).

174

175 *Analyzing nested genome duplications with POInT.*

176 The paramecia studied here and all vertebrates descend from two sequential genome
177 doubling events (hence the “2R” events in vertebrates). As a result, these genomes have an
178 octoploid state relative to the outgroup used. To model such a whole-genome quadruplication
179 (WGQ), we first used a null model (WGQ_n; Figure 2) where losses occur equally from all four
180 subgenomes, but where the loss rate from triplicated and duplicated loci can differ from that seen
181 in quadruplicated loci. To model the two-step formation of a WGQ, we assumed that the first
182 WGD produced an intermediate polyploid genome where all pillars were in state D_{1,3}. Prior to
183 the second WGD, genes could be lost either from subgenome 1 or subgenome 3, such that, when
184 the second WGD occurred, some pillars are quadruplicated, and some are in state D_{1,3}, because
185 they transitioned from D_{1,3} to S₁ and then to D_{1,2} at the second event, and some are similarly in
186 state D_{3,4} (Figure 2).

187 These WGQ models present a challenge because the POInT computation for such an
188 octoploidy with n genomes scales as $O(24^{2n})$. As a result, it is only computationally feasible to
189 analyze two octoploid genomes. However, if the consecutive whole-genome doublings were
190 sufficiently separated in time, POInT can separate them using the two-step model just described.
191 To do so, we compute the posterior probabilities for each subgenome assignment at each pillar.
192 We are interested in pairs of genomic regions that share a high probability of descending from
193 the same original duplicated region created by the first WGD. This origin is marked by those
194 regions having a high probability of belonging either to subgenomes 1 and 2 or to 3 and 4. We
195 thus sought to phase regions from both octoploidies into pairs of regions created by the most
196 recent genome doubling. For the ciliate genomes, we were able to phase the quadruplicated loci
197 into 11,683 pairs of duplicated loci with at least one gene from each genome and where our
198 orthology assignment confidence for assigning extant genes to one of the two subgenomes from
199 the *first* polyploidy event was $\geq 99\%$. Our results are largely consistent with earlier analyses of

200 these genomes, which also suggested that *Paramecium sexaurelia* branched first after the second
 201 event and that the event and its aftermath were marked by RGL and gene conversion (MCGRATH
 202 *et al.* 2014a; MCGRATH *et al.* 2014b). However, those authors argued that the recent WGD was
 203 likely an autopolyploidy because they detected only modest biases in duplicate loss propensities
 204 between syntentic blocks (MCGRATH *et al.* 2014a). POInT's global bias parameter applied to the
 205 larger dataset used here provides significant evidence for biased fractionation; it appears
 206 therefore the recent *Paramecium* event may have in fact been an allopolyploidy. For the
 207 vertebrate 2R events, a model that attempts to phase the 2R duplicates fit the data no better than
 208 did the null model ($P=0.1$, likelihood ratio test with 1 *d.f.*) and so no further phasing was
 209 attempted.

210

211 *Biased fractionation and convergent losses*

212 The full WGD_{bfc-nb} model used for our main analyses includes convergent loss states C₁
 213 and C₂. When we fit a model (WGD_{bfc}) that allows a fractionation bias to also exist between
 214 these two states, we find that that model fits the data no better than the (unbiased) WGD_{bfc-nb}
 215 model (Supplemental Table 5). Hence, the ϵ parameters in Figure 2 only reflect the degree of
 216 bias observed for pillars passing directly from state U to S₁ or S₂. However, the conclusion of the
 217 presence of biased fractionation in these genomes is still strongly supported when models
 218 without convergent losses are used (Supplemental Table 5), even if, in some cases, the ϵ
 219 estimates are somewhat higher for those models.

220

221 *POInT and topological inference.*

222 For the legume WGD, the grass ρ event, the *Paramecium* tetraploidy, the nematode
 223 triploidy and the salmonid WGD, we used POInT to infer the maximum likelihood phylogeny
 224 under the WGD_{bfc-nb} or WGT_{G3} models and an exhaustive tree search (Supplemental Figure 1).
 225 For the Brassica WGT, we assumed that *B. rapa* and *B. oleracea* were sister taxa and tested all
 226 three rooted topologies consistent with this constraint. The topology for the yeast WGD was
 227 taken from Kurtzman and Robnett (2003), for the TGD from Near *et al.*, (2012) and for At- α
 228 from Huang *et al.*, (2016). The vertebrate 2R topology is trivial.

229 For the salmonid WGD, the inferred topology differs significantly from others that have

230 been published. We therefore fit the full POInT model under the topology published by Crespi
 231 and Fulton (2004). The orthology estimates and model parameters are largely unaffected by this
 232 topology change: the orthology relationships of only 106 (0.7%) pillars with posterior probability
 233 >80% differ when the topology is changed, and 91 of these changes simply swap the identities of
 234 the more and less fractionated genomes. The corresponding figures for 95% confidence are 9 and
 235 7 pillars.

236

237 *Orthology inferences and inference of synonymous distances.*

238 Using high confidence orthologs estimated with POInT, we computed the mean
 239 synonymous divergence for every branch for each polyploidy event. The nematode triploidy and
 240 vertebrate 2R events were omitted from this analysis due to their fragmented synteny blocks. For
 241 the tetraploidies, we considered “nearly fully duplicated” pillars: i.e., pillars with at most one
 242 missing gene copy from each of the two gene trees produced by the genome duplication (two
 243 total losses) for all events except the TGD and yeast WGDs, where we allowed two losses from
 244 each subtree (four total losses). For the Brassica hexaploidy, we analyzed only fully triplicated
 245 pillars. At each such pillar, we aligned amino acid sequences for the genes in question with T-
 246 coffee (NOTREDAME *et al.* 2000). We fit the Goldman and Yang codon model of evolution
 247 (GOLDMAN AND YANG 1994) to the corresponding codon-preserving alignments and mirrored
 248 gene trees and extracted the estimated synonymous divergence (K_s) for each branch from this
 249 codon model as described by these authors.

250 With the possible exception of the salmonids and ciliates (ALLENDORF AND THORGAARD
 251 1984; MACREADY *et al.* 1996; BRAASCH AND POSTLETHWAIT 2012), all of the events studied here
 252 are believed to be allopolyploids (THOMAS *et al.* 2006; SCHNABLE *et al.* 2012; TANG *et al.* 2012;
 253 MARCET-HOUBEN AND GABALDON 2015; CONANT 2020; SCHOONMAKER *et al.* 2020). For a
 254 given pillar in set of allopolyploid taxa, the mean synonymous divergence observed along this
 255 root branch ($\overline{K_s^R}$; Figure 3) should represent the sum of the pre-polyploidy divergence of the
 256 diploid progenitors as well as the divergence that occurred after the formation of polyploid but
 257 before the first speciation event among the polyploid taxa. However, recombination events
 258 could, through genetic drift, result in the replacement of alleles from one of the progenitors with
 259 those from the other (WOLFE 2001). These recombinations, or homoeologous exchanges (HE;
 260 GAETA AND CHRIS PIRES 2010) are reasonably common in neopolyploid plants (DOYLE *et al.*

261 2008; CHALHOUB *et al.* 2014; ZHANG *et al.* 2020), but it is not clear whether they are frequent
262 enough to affect the divergence seen along these root branches. We extracted the coding
263 sequences for each pillar that had every homoeologous gene preserved. Post-polyploidy
264 homoeolog displacement (GAUT AND DOEBLEY 1997; WOLFE 2001) will erase the divergence
265 between the progenitor genomes, leaving only the post-displacement divergence to be observed.
266 In such a case, we might expect to observe two modes in synonymous divergence, a larger value
267 for homoeologs that did not experience displacement and a smaller one (lacking the progenitor
268 divergence) for homoeologs that did. To test this hypothesis, we fit the set of estimated
269 synonymous divergences (K_s) along the root branches to either one or two log-normal
270 distributions using the R package *mclust* (SCRUCCA *et al.* 2016) with the best-fit model (i.e., one
271 or two distributions) chosen with the Bayesian information criterion (BIC; SCHWARZ 1978).
272 Values of K_s less than 5×10^{-3} or greater than 2.0 were omitted from these analyses as
273 representing either no synonymous divergence or saturated synonymous divergence,
274 respectively. When two distributions were fit, a “weighting” p reflecting the mixing proportion
275 of each component was also estimated. For a few root branches, a bimodal distribution is
276 preferred. However, in most cases this bimodality is not consistent across different collections of
277 pillars and, even when it is, the proportion of pillars belonging to one of the “modes” is generally
278 very small (Supplemental Table 6). We hence see little suggestion of HE in these data.

279
280 *Filtering for extreme instances of gene conversion.*

281 Because gene conversion among homoeologs (as seen in yeasts; EVANGELISTI AND
282 CONANT 2010; SCIENSKI *et al.* 2015) could confound our K_s estimates, we sought to filter out
283 pillars that showed strong evidence of having experienced it. We created “gene conversion gene
284 trees” for each pillar where each homoeologous gene was forced to be sister to its paralog(s).
285 Any pillars where the likelihood of the sequence alignment under these gene conversion trees
286 was higher than that seen in the mirrored species trees was omitted from our estimates of
287 synonymous divergence (Supplemental Figure 2).

288
289 *Comparing duplicate loss rates to estimated synonymous divergence.*

290 Using the K_s inferences made above for each branch, we compared POInT’s maximum
291 likelihood estimate (MLE) of the rate of homoeolog loss (i.e., its estimated branch length, αt in

292 Supplemental Figure 1) to each branch's mean synonymous divergence, $\overline{K_s}$, to see if the number
293 of losses on any particular branch was unusually large or small. Previous studies that used gene
294 tree approaches to inferring loss rates (TILEY *et al.* 2016) are not comparable to the results here
295 because, unlike POInT, they do not account for the fact that the *observed* homoeolog loss rate
296 necessarily declines in time because progressively fewer homoeolog pairs remain to be lost.
297 Similarly, prior parsimony-based analyses do not include the uncertainty inherent in estimating
298 loss timing, which we account for using POInT's explicit phylogenetic models (MCGRATH *et al.*
299 2014a). Estimating confidence intervals for these ratios of $\alpha \cdot t / \overline{K_s}$ is challenging. We treated
300 the numerators and denominators of these ratios as being normally distributed and independent
301 random variables. The maximum likelihood estimates (MLEs) of αt in the numerators should
302 have asymptotically normal distributions with means that are equal to the true parameter values.
303 The variances of these normal distributions were approximated by evaluating the inverse of the
304 observed Fisher information (i.e., the Hessian of the negative log-likelihood; see KENDALL AND
305 STUART 1973). We estimated the observed Fisher information values via a single-dimension
306 finite difference approximation that ignored covariances between the αt parameter and other
307 parameters.

308 For each branch of the phylogeny, the K_s estimates that are in the denominator of the
309 ratio $\alpha \cdot t / \overline{K_s}$ are obtained via a sample mean of the K_s estimates from the sequences of
310 individual pillars (i.e., $\overline{K_s}$). Due to the Central Limit Theorem, this sample mean should be
311 approximately normally distributed with mean equal to the true parameter value and with
312 variance being approximately the sample variance among individual K_s estimates divided by the
313 number of individual K_s estimates.

314 To infer confidence intervals for the ratio of $\alpha \cdot t / \overline{K_s}$ on each branch, we independently
315 sampled from the two aforementioned normal distributions that are used to approximate the
316 uncertainty of αt and $\overline{K_s}$ estimates in the ratio. For each branch, we calculated the ratio of these
317 sampled values for 1000 pairs of randomly drawn values. We then sorted the resulting ratios and
318 set 95% confidence intervals by finding the ratio value that defined the lower and upper 2.5% of
319 the sorted values.

320 Because the inclusion of fixation in our loss models can give rise to long tip branches
321 (effectively the model suggests that all surviving duplicates in some genomes are now fixed), we
322 present data using a model with convergent losses and biased fractionation but no fixation

323 (WGD_{bc-nb}). However, our results are very similar when using the full WGD_{bfc-nb} model
324 (Supplemental Figure 3).

325
326 *Potential biases in estimating the rate of early duplicate losses.*

327
328 One might object that this signal of rapid early duplicate losses might instead be due to
329 genes being missing from one of the allopolyploid progenitors, meaning that the duplicate pair in
330 question never formed. In this case, the estimates of loss rates along the root branch might be
331 inflated. *A priori*, this idea appears unlikely because of the nature of the genes selected for
332 analysis with POInT. Our inference pipeline requires that each pillar be mapped to a single-copy
333 gene in an outgroup genome (Supplemental Table 1) and that every polyploid genome possess at
334 least one copy of that gene. Hence, the genes analyzed are on average very well conserved over
335 the tree, making a large number of losses of such genes in a progenitor unlikely. Furthermore, in
336 the special case of a hexaploidy, we can actually use POInT to estimate the proportion of genes
337 missing from at least one progenitor genome. Specifically, for the *Brassica* hexaploidy, we
338 showed that the proportion of pillars where a gene was missing from the last-arriving progenitor
339 subgenome (termed LF, or “least fractionated”) was only ~0.3% (HAO *et al.* 2021). Finally, we
340 can also explore the hypothesis of a large number of missing progenitor genes by looking at the
341 patterns of biased fractionation on the root branch relative to the other branches of the tree. We
342 fit a model where the biased fractionation parameter ϵ was allowed to differ on the root branch
343 relative to the other branches, using the WGD_{bf} model above to avoid concerns with convergent
344 losses. Losses from the progenitor genomes prior to polyploidy should be balanced, since biased
345 fractionation is driven by forces that appear at the polyploidy event. Hence, under a model of
346 numerous pre-polyploidy losses, the fractionation bias on the root branch should be *lower* (larger
347 ϵ) than on subsequent branches. Instead, in several cases, the level of biased fractionation is
348 actually higher on the root branch (i.e., the inferred value of ϵ is smaller along the root branch for
349 *At- α* and the paramecium and salmonid WGDs; Supplemental Table 5), consistent with our prior
350 observations in yeast (EMERY *et al.* 2018). Given this fact, and because in some cases upwards of
351 50% of the currently fully single-copy genes in these genomes were returned to single copy
352 along this root branch (Supplemental Figure 4), the degree of pre-polyploidy losses that would
353 be required to bias the results in Figure 4 is implausibly high.

354

355 *Comparisons of selective constraint for different classes of polyploid loci*

356 We examined the inferred average selective constraint (K_a/K_s , estimated as described
357 above) for five classes of polyploid loci (i.e., pillars) across the seven WGD events: 1) Pillars
358 that are single copy in all taxa and have a high probability of having returned to single-copy
359 along the root branch, 2) Pillars that are completely single copy but where the genes did not
360 return to single-copy on the root branch (i.e., where alternative copies of the duplicated genes are
361 preserved in different genomes), 3) pillars with duplicates surviving in only a single species, 4)
362 pillars where all but one species maintains the duplication and 5) pillars where all species
363 maintain duplicate copies. Confidence intervals for these mean K_a/K_s estimates were estimated
364 with the approach described above.

365

366 *Identifying reciprocal gene losses (RGLs) between polyploid taxa.*

367 For a pair of single-copy genes from distinct genomes, the probability that these genes
368 represent RGLs is simply the sum of the probabilities of the orthology relationships, estimated
369 with POInT, that place them as paralogs rather than orthologs. We computed, for each pair of
370 extant taxa in each event, the set of RGLs that we could identify with a confidence of $\geq 95\%$
371 (Figure 5A). To avoid spurious inferences, we restricted our identification of RGL pairs to
372 single-copy genes in each genome where either: a) both the gene and the “hole” corresponding to
373 its lost homoeolog were in synteny with genes on either side or b) the single-copy gene in
374 question was the only homolog of the outgroup gene used for the inference of the NCS blocks. In
375 the first case, this filter corresponds to a clear absence of a corresponding homoeolog in the
376 paralogous synteny block, in the second to the absence of a gene that could be the “missing”
377 homoeolog. We then used TBLASTX (ALTSCHUL *et al.* 1997) to search the non-coding regions
378 of each genome for putative homoeologous copies of the inferred RGL gene that were missed in
379 the genome annotations (i.e., the inference of RGL was spurious due to an annotation artifact). In
380 Case “a” above, this search was restricted to the non-coding regions in the “hole” between the
381 neighboring syntenic genes; in Case “b,” we searched the entire genome for the potentially
382 unannotated homoeolog. Only RGL genes with no such matching noncoding regions at an E-
383 value cutoff of $\leq 10^{-10}$ were considered “true” RGLs. These secondary filters were not applied for
384 the yeast WGD because those data were taken from the manually curated Yeast Genome Order
385 Browser (YGOB, BYRNE AND WOLFE 2005).

386 Data on gene knockouts producing lethal phenotypes from zebrafish, *A. thaliana* and
387 bakers' yeast were taken from ZFIN (HOWE *et al.* 2013; CONANT 2020); a set of 510 “embryo-
388 defective” genes identified by Meinke (2020); and Steinmetz *et al.*, (2002), respectively. The
389 proportion of RGLs in these “essential gene” lists was compared to the proportion of all other
390 single-copy genes from the same organism in the list using Fisher's exact test (SOKAL AND
391 ROHLF 1995). For these same three species, we used GeneOntology data (GENE ONTOLOGY
392 CONSORTIUM 2015) and Panther Overrepresentation Tests (Release 20200728; MI *et al.* 2019) to
393 ask if there were terms from the GO-Slim Biological Process, Cellular Compartment or
394 Molecular Function ontologies that differed in their frequency between the RGL genes and other
395 single-copy genes. After FDR correction (BENJAMINI AND HOCHBERG 1995), no such terms were
396 found for any of the three ontologies across any of the three genomes (FDR-corrected P -value >
397 0.05).

398

399 **Results**

400 *Modeling evolution after ten independent polyploidies.*

401 Using POInT, we assembled a set of ~70,000 homoeologous loci produced by ten
402 different polyploidies. For each polyploidy event, we inferred a set of pillars that it created and
403 ordered them so as to maximize the retained synteny among the extant genes, approximating the
404 ancestral order of the single-copy genes just prior to polyploidy (*Methods*). Six of the events are
405 whole genome duplications (WGDs or tetraploidies): At- α in *Arabidopsis thaliana* and its
406 relatives, a WGD found in legumes, the ρ event from grasses, the teleost-specific genome
407 duplication (TGD), and WGDs from salmonids and yeasts. We further analyzed an asexual
408 triploidy in nematodes, a hexaploidy (whole genome triplication; WGT) in cabbages and their
409 relatives (*Brassica* WGT) and two octoploidies: the vertebrate 2R polyploidy and another in the
410 paramecia (Figure 3). Analyzing octoploidies in POInT is computationally expensive. As a
411 result, we modeled the octoploidy among the paramecia as occurring via two sequential genome
412 duplications and then extracted and analyzed only the more recent of these two events for the
413 remainder of our work (*Methods*). This approach failed with the vertebrate 2R event, presumably
414 because the two events are very ancient and closely spaced in time. A visual interface to these
415 data is available from the POInT browser (<http://wgd.statgen.ncsu.edu>).

416

417
 418
 419
 420
 421
 422
 423
 424
 425
 426
 427
 428
 429

For the WGD events, we compared nested models of evolution (Figure 2 and Supplemental Table 3) that describe the process of homoeolog loss after polyploidy: these models differ as to whether they include biased fractionation, homoeolog fixation and convergent homoeolog losses. For all seven tetraploidies, models that allow for homoeolog pairs to be retained as fixed duplicates after polyploidy fit the observed loss data better than models without such an effect ($\gamma \neq 0$; $P < 10^{-10}$; likelihood ratio test or LRT; Figures 2&3). In addition, every event save that in yeast shows strong evidence for biased fractionation ($\epsilon \neq 1$; $P < 10^{-7}$; LRT; Figure 2&3), while all but the *Paramecium* event show a pattern of independent yet convergent losses to the same homoeolog in independent lineages ($\delta \neq 0$; $P < 10^{-10}$; LRT; Figures 2&3). The nematode triploidy and the *Brassica* WGT also share similar patterns of biased fractionation (Figures 2&3 and Supplemental Table 4).

430
 431
 432
 433
 434
 435
 436
 437
 438
 439
 440
 441
 442

The fact that these events are of widely differing ages is evident from the different degrees of loss/resolution seen in the extant genomes. The branches of Figure 3 are color-coded by POInT's inferences of the proportion of single-copy genes (i.e., loci where all but one of the homoeologous genes have been lost) present at their beginning and ending. While the yeast WGD is inferred to be nearly "fully" resolved (nearly all homeologous loci have been reduced to single-copy or fixed as duplicates), the tetraploidy in salmonid fishes and the nematode triploidy show proportionally few single-copy genes. The nematode triploidy differs from the remaining events in that these animals are asexual triploids and are likely under a different selective regime in their gene losses, (SCHOONMAKER *et al.* 2020). The continued occurrence of meiotic chromosome pairings of homoeologous chromosomes created by the salmonid event may have reduced the rate of homoeolog loss in those genomes (ALLENDORF *et al.* 2015).

443
 444

Many events show rapid homoeolog loss immediately after polyploidy.

445
 446
 447

Loss of duplicate genes immediately after polyploidy can be rapid (SCANNELL *et al.* 2006; SCANNELL *et al.* 2007), and at least two non-exclusive hypotheses exist as to why. The first is that genetic drift should eliminate truly redundant gene copies quickly (LI 1980; LYNCH

448 AND CONERY 2000). The second is the potential for “selected” duplicate losses, an idea suggested
449 by the observation of gene families found to be persistently returned to single-copy after
450 independent polyploidies (PATERSON *et al.* 2006). Such losses might occur if the increases in
451 gene copy number after polyploidy induce disadvantageous dosage conflicts for these genes,
452 such that natural selection acts to remove the homoeologous copies in question (EDGER AND
453 PIRES 2009; DE SMET *et al.* 2013).

454 To study the pattern of early losses, we examined the divergence that occurred
455 immediately after the polyploidy event and prior to any speciation events. In the context of a
456 gene tree for a pair of homoeologous genes produced by a WGD, this period corresponds to the
457 internal branch of the gene tree separating that pair of homoeologs. For a WGT, the situation is
458 analogous except that there are three such branches separating the three homoeologous copies.
459 For simplicity, we refer to these branch(es) as the “root” (purple in Figure 3B). For all branches
460 in each event, we obtained a rough estimate of the time encompassed by that branch by using the
461 mean number of synonymous substitutions per synonymous site (\overline{K}_s) across many homoeologous
462 genes as a neutral clock (*Methods*). The rate of homoeolog loss for each branch is given by
463 POInT’s branch length estimate (αt), computed with its irreversible loss model, such that these
464 branch lengths are scaled based on the number of homoeologous copies at the beginning of that
465 branch (meaning that they are not biased by the fact that later branches have fewer total
466 homoeologs available for loss, *Methods*). The ratio of $\alpha \cdot t / \overline{K}_s$ gives a sense of whether
467 homoeolog losses per time are unusually high or low for a given branch relative to other
468 branches in the same event. For the majority of the polyploidies, we found that the $\alpha \cdot t / \overline{K}_s$ ratio
469 was higher for the root branch than any other branch, consistent with a more rapid loss of
470 homoeologs along this branch (Figure 4). This result is the more striking because the inferred
471 mean K_s value for the root branch (\overline{K}_s^R) should, in the case of an allopolyploidy, also include the
472 pre-polyploidy progenitor divergences. Hence, the \overline{K}_s^R values for these events should be over-
473 estimates, making the $\alpha \cdot t / \overline{K}_s^R$ ratio an underestimate of the relative homoeolog loss rate along
474 the root branch.

475
476
477
478

479 If natural selection were actively favoring the loss of some homoeologous copies
480 immediately after polyploidy, it is possible that the genes involved in those early losses would
481 display a stronger selective constraint than do homoeologous copies lost later in that event's
482 history due to the possibility of dominant negative interactions or expression-linked dosage
483 conflicts (DRUMMOND *et al.* 2006; DE SMET *et al.* 2013; VEITIA *et al.* 2013). We hence compared
484 the average selective constraint, measured as the ratio of nonsynonymous to synonymous
485 substitutions (K_a/K_s), of two types of fully single-copy genes. The first are the single-copy genes
486 whose homoeolog was lost immediately after the polyploidy event along the root branch; the
487 second are the fully single-copy genes where different extant genomes retain homoeologous
488 copies from alternative subgenomes, a situation that requires that the losses have occurred
489 independently after the first speciation event. For most events we observe little difference in
490 constraint between these two groups, while for the Legume WGD the single-copy genes lost later
491 are actually *more* constrained, the opposite of the prediction for selected losses (Supplemental
492 Figure 5).

493

494 *Extensive reciprocal gene loss between pairs of polyploid taxa.*

495 Following Scannell and colleagues (2006; 2007), we searched for post-polyploidy
496 reciprocal gene losses (RGL). We omitted the vertebrate 2R and nematode triploidy from this
497 analysis due to the fragmented nature of the genomes used. With the exception of three closely
498 related yeast species in the *Saccharomyces* genus, every pair of genomes in our remaining eight
499 polyploidies were separated by at least 4 RGLs (this minimal number was seen in the platyfish,
500 tilapia and medaka clade of the TGD; Figure 5C), with the number rising to over a thousand for a
501 few of the yeast taxa pairs. These conclusions are also robust to the confidence cutoffs used to
502 infer the RGLs (Supplemental Figure 6). Our results are in accord with previous work in yeasts
503 and grasses (SCANNELL *et al.* 2006; SCANNELL *et al.* 2007; SCHNABLE *et al.* 2012), and there
504 appears to be a relatively direct relationship between the synonymous divergence of a pair of
505 taxa (a proxy for divergence time) and the number of RGLs separating them (Figure 4A and B).
506 Such a relationship would be expected if both RGLs and synonymous substitutions were
507 accumulating through neutral evolutionary processes (Figure 5A). However, the proportionality
508 between synonymous substitutions and RGLs differs between polyploidy events, with the yeast
509 WGD showing more RGLs per unit K_s than the other events. When we compared the genes

510 involved in reciprocal losses in zebrafish, *A. thaliana* and bakers' yeast to other single-copy
511 genes, there were no significant functional differences between these two sets, again as one
512 would expect were RGL a neutral process (*Methods*).

513
514 The evolutionary importance of RGLs can be assessed by the biological role of the genes
515 that experienced it. For instance, were only “non-essential” genes to experience RGL, then it
516 might not present significant barriers to hybridization. On the other hand, two populations
517 separated by a single RGL for an essential gene would form diploid hybrids whose gametes
518 would lack the gene in question 25% of the time. We can use experimental data on gene
519 essentiality from bakers' yeast, *A. thaliana* and zebrafish (*Methods*) to ask whether the
520 proportion of RGLs that include an essential gene differs from the overall proportion of essential
521 single-copy genes. For the At- α and TGD events, the proportion of RGLs where the surviving
522 gene in *A. thaliana* or zebrafish is essential does not differ from the proportion of other single-
523 copy genes that are essential (Supplemental Table 7). Curiously, the RGLs found when
524 comparing bakers' yeast to some of its nearer relatives are actually *more* likely to be essential
525 than other single-copy genes (Supplemental Table 7). This overrepresentation is likely due to
526 the fact that the duplicate losses that occurred prior to the first speciation event were actually
527 underrepresented in essential genes (Supplemental Table 8). As a result, RGLs, which must have
528 occurred *after* the first speciation event (see the yeast clade of Figure 3), would be enriched in
529 essential genes simply because more essential genes survived in duplicate past that first
530 speciation.

531
532 The importance of RGL in driving speciation events among polyploid taxa has been
533 questioned on theoretical grounds, as the appearance of RGLs is subject to the same requirement
534 of reproductive isolation as are the appearances of other genetic incompatibilities among
535 populations (MUIR AND HAHN 2015). This objection has more force for obligately sexual
536 organisms than it does for organisms such as bakers' yeast, where it is estimated that there are
537 1000 mitotic cell divisions for every meiosis and that only about 1% of meioses are out-crosses
538 (TSAI *et al.* 2008). Indeed, Figure 5 suggests that RGL may occur more frequently in yeasts (and
539 potentially in some plants, which may also reproduce asexually) than in the teleost fishes and
540 particularly the salmonids.

541 Even if RGL does not drive speciation, it still represents a barrier to diploid hybrids: most
542 of the taxa pairs for which essentiality data are available are separated from each other by at least
543 one RGL for an essential gene, the exceptions being some of the closest relatives of *A. thaliana*,
544 zebrafish and bakers' yeast studied (Supplemental Table 7). This observation is consistent with
545 studies of the relative frequency of diploid and polyploid hybridizations in flowering plants. In
546 these lineages, it is rare to find successful diploid hybrids involving distantly related parental
547 species (where RGLs could be common). However, allopolyploid hybrids appear to form at
548 roughly the same rate across a much larger range of divergence times (BUGGS *et al.* 2009b). A
549 potential explanation for the frequency of recurrent polyploidy is therefore simply that a new
550 allopolyploidy can allow paleopolyploids to again enjoy the benefits of hybridization (such as
551 hybrid vigor and heterosis; BIRCHLER *et al.* 2006; CHEN 2010) in the face of their isolation due to
552 RGL.

553

554 Discussion

555 There are a surprising number of similarities seen in the manner of polyploidy resolution
556 across these independent polyploidies. Biased fractionation and other patterns in the homoeolog
557 losses are similar across many events: reciprocal gene losses are also present for most pairs of
558 polyploid taxa. The rate of homoeolog loss immediately after polyploidy is very high for many,
559 but not all, events (Figure 4).

560 Moreover, the differences in evolutionary patterns we do see are often in keeping with
561 what we know about the history of the events themselves. For instance, the salmonid WGD is
562 marked by continuing pairing of homoeologous chromosomes in meiosis (ALLENDORF *et al.*
563 2015). These pairings appear to limit the number of homoeolog losses and, for this event, loss
564 rates at the phylogeny tips and root are similar (per unit K_s). The grass ρ and yeast events have
565 loss rates that are roughly similar (again per unit K_s) across time, a fact for which we currently
566 do not have an operating hypothesis.

567 For the events that do show rapid losses along the root branch, which of the two
568 hypotheses mentioned, drift or selected losses, seems to better explain our data? The homoeologs
569 lost along the root are not more selectively constrained than other purely single-copy genes
570 known to have been lost later (Supplemental Figure 5). This fact probably speaks against any
571 very large number of selected losses. The single-copy genes as a whole are also generally

572 somewhat less selectively constrained than are genes with surviving homoeologs (Supplemental
573 Figure 5). Moreover, there is a clear pattern in most events whereby most of the fully single-copy
574 genes that exist today are predicted to have lost their homoeologous partner along the root
575 branch (Supplemental Figure 4). The yeast, nematode, and Paramecium events may violate this
576 pattern because the nematode event is an asexual triploidy while the other two involve lineages
577 that have significant rates of asexual reproduction. In such cases, restoring proper meiotic pairing
578 is less necessary than in taxa with primarily sexual reproduction. As a result, we expect that
579 asexually reproducing lineages could more easily form viable new species immediately after
580 polyploidy, meaning that the post-polyploid “lag” in speciation might be less evident (SCHRANZ
581 *et al.* 2012). As a preliminary hypothesis, we therefore propose that, for most polyploidies in
582 animals and plants, the majority of the purely neutral homoeolog losses occur prior to extensive
583 species divergence in the polyploid clade. A natural extension to this proposal would be that the
584 post-polyploidy lag represents this earlier period of neutral homoeolog loss, though the question
585 of why speciation events might be rare during such a period is still to be answered. A further
586 implication would be that later losses (including RGLs) would have occurred in homoeologous
587 pairs that were initially preserved to maintain dosage balance. They are then only lost when later
588 mutations, such as expression changes, release this dosage constraint and allow the loss of one of
589 the copies (BIRCHLER *et al.* 2005; CONANT *et al.* 2014). The higher selective constraint of genes
590 with surviving homoeologs is arguably also consistent with this hypothesis.

591 While the best-studied ancient polyploidy is in bakers’ yeast, it is atypical in a number of
592 respects. Biased fractionation is much less evident here (EMERY *et al.* 2018), losses are not
593 heavily biased toward the earliest phases of the polyploidy (Figure 4) and RGL is much more
594 prevalent. As mentioned above, one major source of these differences is likely the relative timing
595 of the post-polyploidy speciations: the yeasts had almost no lag between their polyploidy event
596 and the first observed speciation in our dataset (Supplemental Figure 4; SCHRANZ *et al.* 2012).

597 Other questions remain unanswered. The relative formation rates of allo- and
598 autopolyploids are uncertain. While recent polyploids appear to be approximately equally
599 divided between the two (BARKER *et al.* 2016), the potential selective advantages of being an
600 allopolyploid, and hence a hybrid (ALIX *et al.* 2017; BLANC-MATHIEU *et al.* 2017), could result
601 in a strong skew towards allopolyploids among the rare polyploidies that survive to become the
602 ancient events of the kind studied here (BARKER *et al.* 2016). The results here are consistent with

603 this hypothesis, but our sample of events is potentially biased by the available genome
604 sequences. Across all of the events, we find that the ubiquity of homoeolog fixation and (except
605 in paramecia) convergent homoeolog losses both speak to a common selective environment
606 acting to maintain certain homoeologs after all of these events. The most obvious candidate for
607 such a selective force is again the dosage balance hypothesis: it argues that highly interacting
608 genes tend to remain in multiple copies post-polyploidy to preserve the stoichiometry of those
609 interactions (BIRCHLER *et al.* 2005; BIRCHLER AND VEITIA 2012; TASHIGHIAN *et al.* 2017).
610 Whatever the role of RGL in speciation, it is clear that all of these polyploid organisms possess a
611 degree of isolation due to it. The role of RGL in recurrent polyploidy is hence an important topic
612 for future research. Biology has a history of viewing “rules” as being more honored in the
613 breach, but the commonalities in post-polyploidy genome evolution across wide taxonomic
614 distances are both interesting in their own right and for the insight they give on other aspects of
615 biology (PIRES AND CONANT 2016).

616

617

618

619 **Data availability:**

620 All underlying data are available from the POInT browser (wgd.statgen.ncsu.edu) and from
621 figshare (DOI: <https://doi.org/10.6084/m9.figshare.12750992.v4>); the POInT package (v1.55) is
622 available from GitHub (<https://github.com/gconant0/POInT>)

623

624 **Acknowledgements:**

625 We would like to thank K. Wolfe for helpful comments and K. Byrne for help with the
626 YGOB datasets.

627

628 **Funding:**

629 YH, JCP and GCC were supported by National Science Foundation grant NSF-IOS-1339156.
630 EL was supported by NSF-IOS-1339156 and NSF-IOS-1849708. JLT was supported by NSF-
631 DEB-1754142 and by National Institutes of Health grant NIH-R01-GM118508.

632

633 **Competing interests:** The authors declare that they have no competing interests.

634 **References:**

635

636 Alix, K., P. R. Gérard, T. Schwarzacher and J. Heslop-Harrison, 2017 Polyploidy and interspecific
637 hybridization: partners for adaptation, speciation and evolution in plants. *Annals of botany* 120:
638 183-194.

639 Allendorf, F. W., S. Bassham, W. A. Cresko, M. T. Limborg, L. W. Seeb *et al.*, 2015 Effects of
640 crossovers between homeologs on inheritance and population genomics in polyploid-derived
641 salmonid fishes. *Journal of Heredity* 106: 217-227.

642 Allendorf, F. W., and G. H. Thorgaard, 1984 Tetraploidy and the evolution of salmonid fishes, pp. 1-53 in
643 *Evolutionary genetics of fishes*. Springer.

644 Altschul, S. F., T. L. Madden, A. A. Schaffer, J. H. Zhang, Z. Zhang *et al.*, 1997 Gapped Blast and Psi-
645 Blast : A new-generation of protein database search programs. *Nucleic Acids Research* 25: 3389-
646 3402.

647 Barker, M. S., N. Arrigo, A. E. Baniaga, Z. Li and D. A. Levin, 2016 On the relative abundance of
648 autopolyploids and allopolyploids. *New Phytol* 210: 391-398.

649 Benjamini, Y., and Y. Hochberg, 1995 Controlling the false discovery rate: A practical and powerful
650 approach to multiple testing. *Journal of the Royal Statistical Society, Series B (Methodological)*
651 57: 289-300.

652 Birchler, J. A., N. C. Riddle, D. L. Auger and R. A. Veitia, 2005 Dosage balance in gene regulation:
653 biological implications. *Trends Genet* 21: 219-226.

654 Birchler, J. A., and R. A. Veitia, 2012 Gene balance hypothesis: connecting issues of dosage sensitivity
655 across biological disciplines. *Proc Natl Acad Sci U S A* 109: 14746-14753.

656 Birchler, J. A., H. Yao and S. Chudalayandi, 2006 Unraveling the genetic basis of hybrid vigor.
657 *Proceedings of the National Academy of Sciences* 103: 12957-12958.

658 Blanc, G., and K. H. Wolfe, 2004 Functional divergence of duplicated genes formed by polyploidy during
659 Arabidopsis evolution. *Plant Cell* 16: 1679-1691.

660 Blanc-Mathieu, R., L. Perfus-Barbeoch, J.-M. Aury, M. Da Rocha, J. Gouzy *et al.*, 2017 Hybridization
661 and polyploidy enable genomic plasticity without sex in the most devastating plant-parasitic
662 nematodes. *PLoS Genetics* 13: e1006777.

663 Braasch, I., and J. H. Postlethwait, 2012 Polyploidy in fish and the teleost genome duplication, pp. 341-
664 383 in *Polyploidy and genome evolution*. Springer.

665 Buggs, R., A. Doust, J. Tate, J. Koh, K. Soltis *et al.*, 2009a Gene loss and silencing in *Tragopogon*
666 *miscellus* (Asteraceae): comparison of natural and synthetic allotetraploids. *Heredity* 103: 73-81.

667 Buggs, R. J., P. S. Soltis and D. E. Soltis, 2009b Does hybridization between divergent progenitors drive
668 whole-genome duplication? *Molecular Ecology* 18: 3334-3339.

669 Byrne, K. P., and K. H. Wolfe, 2005 The Yeast Gene Order Browser: Combining curated homology and
670 syntenic context reveals gene fate in polyploid species. *Genome Research* 15: 1456-1461.

671 Chalhoub, B., F. Denoeud, S. Liu, I. A. Parkin, H. Tang *et al.*, 2014 Early allopolyploid evolution in the
672 post-Neolithic *Brassica napus* oilseed genome. *science* 345: 950-953.

673 Chen, Z. J., 2010 Molecular mechanisms of polyploidy and hybrid vigor. *Trends in plant science* 15: 57-
674 71.

675 Clausen, R., and T. Goodspeed, 1925 Interspecific hybridization in *Nicotiana*. II. A tetraploid *glutinosa-*
676 *tabacum* hybrid, an experimental verification of Winge's hypothesis. *Genetics* 10: 278.

- 677 Conant, G. C., 2014 Comparative genomics as a time machine: How relative gene dosage and metabolic
678 requirements shaped the time-dependent resolution of yeast polyploidy. *Molecular Biology and*
679 *Evolution* 31: 3184-3193.
- 680 Conant, G. C., 2020 The lasting after-effects of an ancient polyploidy on the genomes of teleosts. *Plos*
681 *one* 15: e0231356.
- 682 Conant, G. C., J. A. Birchler and J. C. Pires, 2014 Dosage, duplication, and diploidization: clarifying the
683 interplay of multiple models for duplicate gene evolution over time. *Current opinion in plant*
684 *biology* 19: 91-98.
- 685 Conant, G. C., and A. Wagner, 2002 GenomeHistory: A software tool and its application to fully
686 sequenced genomes. *Nucleic Acids Research* 30: 3378-3386.
- 687 Conant, G. C., and K. H. Wolfe, 2008 Probabilistic cross-species inference of orthologous genomic
688 regions created by whole-genome duplication in yeast. *Genetics* 179: 1681-1692.
- 689 Crespi, B. J., and M. J. Fulton, 2004 Molecular systematics of Salmonidae: combined nuclear data yields
690 a robust phylogeny. *Molecular phylogenetics and evolution* 31: 658-679.
- 691 De Smet, R., K. L. Adams, K. Vandepoele, M. C. Van Montagu, S. Maere *et al.*, 2013 Convergent gene
692 loss following gene and genome duplications creates single-copy families in flowering plants.
693 *Proceedings of the National Academy of Sciences, U.S.A.* 110: 2898-2903.
- 694 Doyle, J. J., L. E. Flagel, A. H. Paterson, R. A. Rapp, D. E. Soltis *et al.*, 2008 Evolutionary genetics of
695 genome merger and doubling in plants. *Annual review of genetics* 42: 443-461.
- 696 Drummond, D. A., A. Raval and C. O. Wilke, 2006 A single determinant dominates the rate of yeast
697 protein evolution. *Molecular Biology and Evolution* 23: 327-337.
- 698 Edger, P. P., and J. C. Pires, 2009 Gene and genome duplications: the impact of dosage-sensitivity on the
699 fate of nuclear genes. *Chromosome research : an international journal on the molecular,*
700 *supramolecular and evolutionary aspects of chromosome biology* 17: 699-717.
- 701 Emery, M., M. M. S. Willis, Y. Hao, K. Barry, K. Oakgrove *et al.*, 2018 Preferential retention of genes
702 from one parental genome after polyploidy illustrates the nature and scope of the genomic
703 conflicts induced by hybridization. *PLoS Genetics* 14: e1007267em.
- 704 Evangelisti, A. M., and G. C. Conant, 2010 Nonrandom survival of gene conversions among yeast
705 ribosomal proteins duplicated through genome doubling. *Genome Biology and Evolution* 2: 826-
706 834.
- 707 Fawcett, J. A., S. Maere and Y. Van de Peer, 2009 Plants with double genomes might have had a better
708 chance to survive the Cretaceous-Tertiary extinction event. *Proc Natl Acad Sci U S A* 106: 5737-
709 5742.
- 710 Felsenstein, J., 1985 Phylogenies and the comparative method. *American Naturalist*: 1-15.
- 711 Freeling, M., 2009 Bias in plant gene content following different sorts of duplication: tandem, whole-
712 genome, segmental, or by transposition. *Annual Review of Plant Biology* 60: 433-453.
- 713 Gaeta, R. T., and J. Chris Pires, 2010 Homoeologous recombination in allopolyploids: the polyploid
714 ratchet. *New Phytologist* 186: 18-28.
- 715 Garsmeur, O., J. C. Schnable, A. Almeida, C. Jourda, A. D'Hont *et al.*, 2013 Two Evolutionarily Distinct
716 Classes of Paleopolyploidy. *Molecular Biology and Evolution* 31: 448-454.
- 717 Gaut, B. S., and J. F. Doebley, 1997 DNA sequence evidence for the segmental allotetraploid origin of
718 maize. *Proceedings of the National Academy of Sciences* 94: 6809-6814.
- 719 Gene Ontology Consortium, 2015 Gene ontology consortium: going forward. *Nucleic acids research* 43:
720 D1049-D1056.

- 721 Goldman, N., and Z. Yang, 1994 A codon-based model of nucleotide substitution for protein-coding
722 DNA sequences. *Molecular Biology and Evolution* 11: 725-736.
- 723 Hao, Y., M. E. Mabry, P. Edger, M. Freeling, C. Zheng *et al.*, 2021 The contributions of the allopolyploid
724 parents of the mesopolyploid Brassiceae are evolutionarily distinct but functionally compatible.
725 *Genome Research* 31: 799-810.
- 726 Howe, D. G., Y. M. Bradford, T. Conlin, A. E. Eagle, D. Fashena *et al.*, 2013 ZFIN, the Zebrafish Model
727 Organism Database: increased support for mutants and transgenics. *Nucleic Acids Research* 41:
728 D854-860.
- 729 Huang, C.-H., R. Sun, Y. Hu, L. Zeng, N. Zhang *et al.*, 2016 Resolution of Brassicaceae phylogeny using
730 nuclear genes uncovers nested radiations and supports convergent morphological evolution.
731 *Molecular biology and evolution* 33: 394-412.
- 732 Kendall, M., and A. Stuart, 1973 *The advanced theory of statistics*. Charles Griffen, London.
- 733 Kirkpatrick, S., C. D. J. Gelatt and M. P. Vecchi, 1983 Optimization by simulated annealing. *Science*
734 220: 671-680.
- 735 Kurtzman, C. P., and C. J. Robnett, 2003 Phylogenetic relationships among yeasts of the '*Saccharomyces*
736 complex' determined from multigene sequence analyses. *FEMS Yeast Research* 3: 417-432.
- 737 Kuwada, Y., 1911 Maiosis in the Pollen Mother Cells of *Zea Mays* L.(With Plate V.). *植物学雑誌* 25:
738 163-181.
- 739 Lander, E. S., and P. Green, 1987 Construction of multilocus genetic linkage maps in humans.
740 *Proceedings of the National Academy of Sciences, U.S.A.* 84: 2363-2367.
- 741 Li, W.-H., 1980 Rate of gene silencing at duplicate loci: A theoretical study and interpretation of data
742 from tetraploid fish. *Genetics* 95: 237-258.
- 743 Lynch, M., and J. S. Conery, 2000 The evolutionary fate and consequences of duplicate genes. *Science*
744 290: 1151-1155.
- 745 Maclean, C. J., and D. Greig, 2011 Reciprocal gene loss following experimental whole-genome
746 duplication causes reproductive isolation in yeast. *Evolution: International Journal of Organic*
747 *Evolution* 65: 932-945.
- 748 Macready, W. G., A. G. Siapas and S. A. Kauffman, 1996 Criticality and Parallelism in Combinatorial
749 Optimization. *Science* 271: 56-59.
- 750 Maere, S., S. De Bodt, J. Raes, T. Casneuf, M. Van Montagu *et al.*, 2005 Modeling gene and genome
751 duplications in eukaryotes. *Proceedings of the National Academy of Sciences, U.S.A.* 102: 5454-
752 5459.
- 753 Marcet-Houben, M., and T. Gabaldon, 2015 Beyond the Whole-Genome Duplication: Phylogenetic
754 Evidence for an Ancient Interspecies Hybridization in the Baker's Yeast Lineage. *PLoS biology*
755 13: e1002220.
- 756 Mayrose, I., S. H. Zhan, C. J. Rothfels, K. Magnuson-Ford, M. S. Barker *et al.*, 2011 Recently formed
757 polyploid plants diversify at lower rates. *Science* 333: 1257.
- 758 McGrath, C. L., J.-F. Gout, T. G. Doak, A. Yanagi and M. Lynch, 2014a Insights into three whole-
759 genome duplications gleaned from the *Paramecium caudatum* genome sequence. *Genetics* 197:
760 1417-1428.
- 761 McGrath, C. L., J.-F. Gout, P. Johri, T. G. Doak and M. Lynch, 2014b Differential retention and
762 divergent resolution of duplicate genes following whole-genome duplication. *Genome research*
763 24: 1665-1675.

- 764 Meinke, D. W., 2020 Genome-wide identification of EMBRYO-DEFECTIVE (EMB) genes required for
765 growth and development in Arabidopsis. *New Phytologist* 226: 306-325.
- 766 Mi, H., A. Muruganujan, D. Ebert, X. Huang and P. D. Thomas, 2019 PANTHER version 14: more
767 genomes, a new PANTHER GO-slim and improvements in enrichment analysis tools. *Nucleic
768 Acids Res* 47: D419-D426.
- 769 Mizuta, Y., Y. Harushima and N. Kurata, 2010 Rice pollen hybrid incompatibility caused by reciprocal
770 gene loss of duplicated genes. *Proceedings of the National Academy of Sciences* 107: 20417-
771 20422.
- 772 Muir, C. D., and M. W. Hahn, 2015 The limited contribution of reciprocal gene loss to increased
773 speciation rates following whole-genome duplication. *The American Naturalist* 185: 70-86.
- 774 Near, T. J., R. I. Eytan, A. Dornburg, K. L. Kuhn, J. A. Moore *et al.*, 2012 Resolution of ray-finned fish
775 phylogeny and timing of diversification. *Proceedings of the National Academy of Sciences,
776 U.S.A.* 109: 13698-13703.
- 777 Notredame, C., D. G. Higgins and J. Heringa, 2000 T-Coffee: A novel method for fast and accurate
778 multiple sequence alignment. *Journal of Molecular Biology* 302: 205-217.
- 779 Ohno, S., 1970 *Evolution by gene duplication*. Springer, New York.
- 780 Orr, H. A., 1996 Dobzhansky, Bateson, and the genetics of speciation. *Genetics* 144: 1331.
- 781 Paterson, A. H., B. A. Chapman, J. C. Kissinger, J. E. Bowers, F. A. Feltus *et al.*, 2006 Many gene and
782 domain families have convergent fates following independent whole-genome duplication events
783 in Arabidopsis, Oryza, Saccharomyces and Tetraodon. *Trends in genetics* 22: 597-602.
- 784 Pires, J. C., and G. C. Conant, 2016 Robust Yet Fragile: Expression Noise, Protein Misfolding and Gene
785 Dosage in the Evolution of Genomes. *Annual Review of Genetics* 50: 113–131.
- 786 Sankoff, D., and M. Blanchette, 1998 Multiple genome rearrangement and breakpoint phylogeny. *Journal
787 of Computational Biology* 5: 555-570.
- 788 Scannell, D. R., K. P. Byrne, J. L. Gordon, S. Wong and K. H. Wolfe, 2006 Multiple rounds of speciation
789 associated with reciprocal gene loss in polyploid yeasts. *Nature* 440: 341-345.
- 790 Scannell, D. R., A. C. Frank, G. C. Conant, K. P. Byrne, M. Woolfit *et al.*, 2007 Independent sorting-out
791 of thousands of duplicated gene pairs in two yeast species descended from a whole-genome
792 duplication. *Proceedings of the National Academy of Sciences, U.S.A.* 104: 8397-8402.
- 793 Schnable, J. C., M. Freeling and E. Lyons, 2012 Genome-wide analysis of syntenic gene deletion in the
794 grasses. *Genome Biol Evol* 4: 265-277.
- 795 Schoonmaker, A., Y. Hao, D. Bird and G. C. Conant, 2020 A single, shared triploidy in three species of
796 parasitic nematodes. *G3: Genes, Genomes, Genetics* 10: 225-233.
- 797 Schranz, M. E., S. Mohammadin and P. P. Edger, 2012 Ancient whole genome duplications, novelty and
798 diversification: the WGD Radiation Lag-Time Model. *Current opinion in plant biology* 15: 147-
799 153.
- 800 Schwarz, G., 1978 Estimating the dimension of a model. *Annals of statistics* 6: 461-464.
- 801 Scienski, K., J. C. Fay and G. C. Conant, 2015 Patterns of Gene Conversion in Duplicated Yeast Histones
802 Suggest Strong Selection on a Coadapted Macromolecular Complex. *Genome Biology and
803 Evolution* 7: 3249-3258.
- 804 Scrucca, L., M. Fop, T. B. Murphy and A. E. Raftery, 2016 mclust 5: clustering, classification and density
805 estimation using Gaussian finite mixture models. *The R journal* 8: 289.
- 806 Seoighe, C., and K. H. Wolfe, 1999 Yeast genome evolution in the post-genome era. *Current Opinion in
807 Microbiology* 2: 548-554.

- 808 Sokal, R. R., and F. J. Rohlf, 1995 *Biometry: 3rd Edition*. W. H. Freeman and Company, New York.
- 809 Soltis, D. E., V. A. Albert, J. Leebens-Mack, C. D. Bell, A. H. Paterson *et al.*, 2009 Polyploidy and
810 angiosperm diversification. *American Journal of Botany* 96: 336-348.
- 811 Soltis, D. E., M. C. Segovia-Salcedo, I. Jordon-Thaden, L. Majure, N. M. Miles *et al.*, 2014a Are
812 polyploids really evolutionary dead-ends (again)? A critical reappraisal of Mayrose *et al.*(2011).
813 *New Phytologist* 202: 1105-1117.
- 814 Soltis, D. E., C. J. Visger and P. S. Soltis, 2014b The polyploidy revolution then... and now: Stebbins
815 revisited. *American journal of botany* 101: 1057-1078.
- 816 Stayton, C. T., 2015 The definition, recognition, and interpretation of convergent evolution, and two new
817 measures for quantifying and assessing the significance of convergence. *Evolution* 69: 2140-
818 2153.
- 819 Stebbins Jr, G. L., 1947 Types of polyploids: their classification and significance, pp. 403-429 in
820 *Advances in genetics*. Elsevier.
- 821 Steinmetz, L. M., C. Scharfe, A. M. Deutschbauer, D. Mokranjac, Z. S. Herman *et al.*, 2002 Systematic
822 screen for human disease genes in yeast. *Nature Genetics* 31: 400-404.
- 823 Tang, H., M. R. Woodhouse, F. Cheng, J. C. Schnable, B. S. Pedersen *et al.*, 2012 Altered patterns of
824 fractionation and exon deletions in *Brassica rapa* support a two-step model of paleohexaploidy.
825 *Genetics* 190: 1563-1574.
- 826 Tasdighian, S., M. Van Bel, Z. Li, Y. Van de Peer, L. Carretero-Paulet *et al.*, 2017 Reciprocally retained
827 genes in the angiosperm lineage show the hallmarks of dosage balance sensitivity. *The Plant Cell*
828 29: 2766-2785.
- 829 Thomas, B. C., B. Pedersen and M. Freeling, 2006 Following tetraploidy in an *Arabidopsis* ancestor,
830 genes were removed preferentially from one homeolog leaving clusters enriched in dose-sensitive
831 genes. *Genome Research* 16: 934-946.
- 832 Tiley, G. P., C. Ané and J. G. Burleigh, 2016 Evaluating and Characterizing Ancient Whole-Genome
833 Duplications in Plants with Gene Count Data. *Genome biology and evolution* 8: 1023-1037.
- 834 Tsai, I. J., D. Bensasson, A. Burt and V. Koufopanou, 2008 Population genomics of the wild yeast
835 *Saccharomyces paradoxus*: Quantifying the life cycle. *Proceedings of the National Academy of*
836 *Sciences, U.S.A.* 105: 4957-4962.
- 837 Van de Peer, Y., T.-L. Ashman, P. S. Soltis and D. E. Soltis, 2021 Polyploidy: an evolutionary and
838 ecological force in stressful times. *The Plant Cell* 33: 11-26.
- 839 Van de Peer, Y., E. Mizrachi and K. Marchal, 2017 The evolutionary significance of polyploidy. *Nature*
840 *Reviews Genetics* 18: 411-424.
- 841 Veitia, R. A., S. Bottani and J. A. Birchler, 2013 Gene dosage effects: nonlinearities, genetic interactions,
842 and dosage compensation. *Trends Genet* 29: 385-393.
- 843 Wagner Jr, W., 1970 Biosystematics and evolutionary noise. *Taxon* 19: 146-151.
- 844 Werth, C. R., and M. D. Windham, 1991 A model for divergent, allopatric speciation of polyploid
845 pteridophytes resulting from silencing of duplicate-gene expression. *The American Naturalist*
846 137: 515-526.
- 847 Wolfe, K. H., 2001 Yesterday's polyploids and the mystery of diploidization. *Nat Rev Genet* 2: 333-341.
- 848 Woodhouse, M. R., J. C. Schnable, B. S. Pedersen, E. Lyons, D. Lisch *et al.*, 2010 Following tetraploidy
849 in maize, a short deletion mechanism removed genes preferentially from one of the two
850 homeologs. *PLoS biology* 8: e1000409.

851 Zhang, Z., X. Gou, H. Xun, Y. Bian, X. Ma *et al.*, 2020 Homoeologous exchanges occur through
852 intragenic recombination generating novel transcripts and proteins in wheat and other polyploids.
853 Proceedings of the National Academy of Sciences.

854

855

856 Figure Legends:

857

858 **Figure 1:** Inferring orthologous chromosome regions between polyploid genomes with POInT. **A)**
859 Cartoon of gene losses and a speciation event after a whole-genome duplication. Immediately after
860 the WGD, all five genes are present in two homoeologous copies. Three homoeologous gene losses
861 occur prior to the split of the two species, one in the less fractionated subgenome (Track “0;” yielding
862 the green gene in the lower window) and two from the more fractionated subgenome (Track “1;”
863 yielding the two blue genes in the upper window). After the speciation event, Genome 1 loses a
864 homoeolog from the more fractionated subgenome and Genome 2 loses one from the less
865 fractionated subgenome, a case of reciprocal gene loss (RGL). **B)** There are $2^n=2^2=4$ potential ways
866 of phasing the chromosomal regions from Genome 1 relative to Genome 2 (i.e., of assigning
867 orthology between the two regions). We identify these 4 states with the subgenome assignment for
868 the top track for each of the two genomes (00→11; red boxes at the right of each diagram). POInT
869 uses a model of homoeolog loss to compute the likelihood of the observed gene presence/absence
870 data at each locus (or “pillar”) for each of these 2^n relationships. These relationships each constitute a
871 hidden state of the HMM implemented by POInT whereas a likelihood of observed gene
872 presence/absence data for a relationship represents an emission probability for the HMM. **C)**
873 Recurrence equation for computing the likelihood of each orthology assignment at pillar i conditional
874 on the data at pillars 0 through $i-1$ (see **B**). For pillar i , we define a vector L^i to be the likelihood of the
875 orthology states, with elements $L_{00}^i, L_{01}^i, L_{10}^i$ and L_{11}^i being POInT’s estimates of the likelihood of each
876 such state based on the gene presence/absence data at that pillar. We then use a transition
877 probability matrix Θ , with each entry representing the probability that pillar i has a particular orthology
878 state conditional upon another orthology state at $i-1$. The probability that the orthology state is
879 maintained between pillars $i-1$ and i is $1-\theta_i$ for each genome (and $(1-\theta_i)^2$ in total); the chance that one
880 genome changes orthology state is $\theta_i(1-\theta_i)$ and the chance that both change is θ_i^2 . Here, $\theta_i=\theta$, a global
881 constant estimated from the data by maximum likelihood, except when synteny is not maintained
882 between pillars, in which case $\theta_i=0.5$ (i.e., adjacent pillars do not inform on each other’s orthology
883 state; *Methods*). To compute a likelihood for the entire data set, POInT implements an HMM forward
884 algorithm that expresses $L^{i|D_i \dots D_0}$, the probabilities of orthology relationships for pillar i and the
885 observed data at pillars 0 through i (denoted $D_i \dots D_0$), in terms of the emission probabilities L^i , the
886 transition probabilities Θ and the probabilities $L^{i-1|D_{i-1} \dots D_0}$ that were already computed for pillar $i-1$.
887 The vector of $L^{i|D_i \dots D_0}$ is then the element-wise vector product (indicated with the “ \odot ”) of $\Theta \cdot$
888 $L^{i-1|D_{i-1} \dots D_0}$ and L^i . This formula can be applied sequentially starting at pillar 0, with the base case
889 $L^{0|D_0} = L^0$. For m pillars, the overall likelihood of the dataset is then the sum of the elements of
890 $L^{m-1|D_{m-1} \dots D_0}$. **D)** POInT employs posterior decoding to infer the orthology relationships at each pillar.
891 Here we illustrate a small region from the most recent *Paramecium* WGD (after phasing from the
892 earlier duplication, see *Methods*), showing the set of orthology relationships inferred by posterior
893 decoding. Genes in adjacent pillars that are also neighbors in an extant genome are shown
894 connected by lines. The number above each pillar is the posterior probability of the inferred orthology
895 relationship. The upper set of three tracks correspond to the less-fractionated parental subgenome,
896 the lower three to the more fractionated one. Gene retained from *only* the less-fractionated genome
897 are colored blue, from *only* the more fractionated one green, and fully retained duplicates are shown
898 in pink. All other patterns of duplicate retention are shown in beige.

899

900 **Figure 2:** Models of polyploidy resolution for three types of events: WGD or whole-genome
901 duplication/tetraploidy, WGT or whole-genome triplication/hexaploidy and WGQ or whole-genome
902 quadruplication/octoploidy). **WGD:** all pillars start in state **U** (**U**ndifferentiated), from which they can

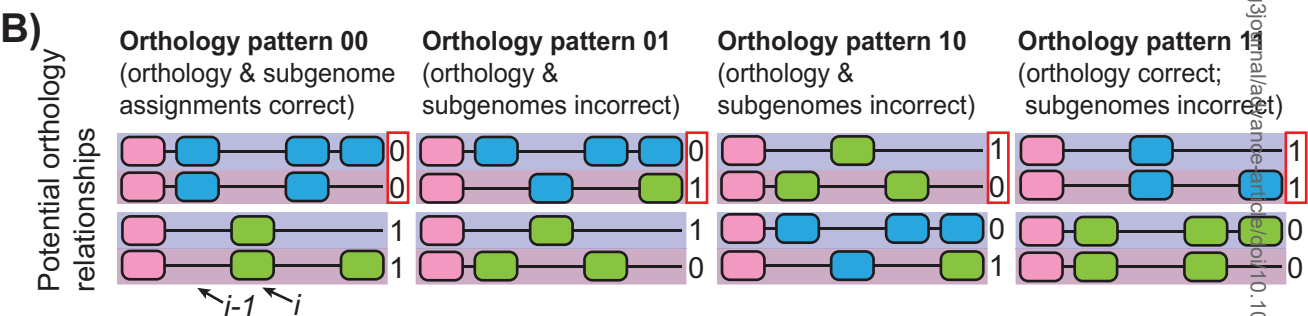
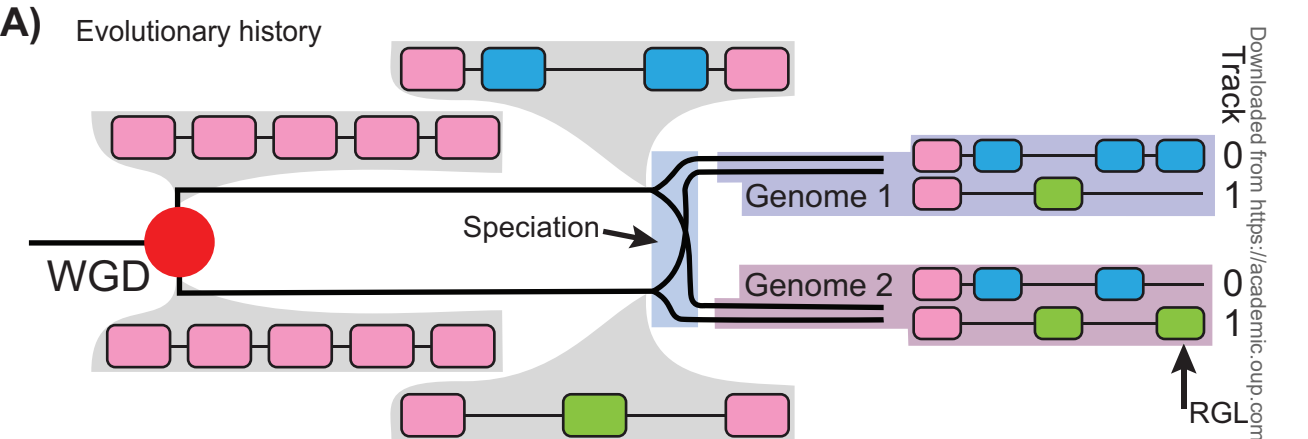
903 transition to either the three other duplicated states, **C**₁ (Converging state 1), **C**₂ (Converging state 2)
 904 and **F** (Fixed) or to the two single-copy states **S**₁ (Single-copy 1) and **S**₂ (Single-copy 2). **C**₁ and **S**₁ are
 905 states where the gene from the less-fractionated parental subgenome will be or are preserved, and
 906 **C**₂ and **S**₂ the corresponding states for the more-fractionated parental subgenome. The null model
 907 has parameters $\gamma=\delta=0$ and $\varepsilon=1.0$. Homoeolog pair fixation is inferred when $\gamma\neq 0$, convergent losses
 908 when $\delta\neq 0$ and biased fractionation when $\varepsilon<1.0$. **WGT**: in the base model all pillars start in state **T**
 909 (**T**riplicated) and transition first to duplicated states (**D**_{x,y}) and hence to the single-copy states (**S**_x).
 910 Genome 1 is assumed to be favored (fewer losses) and the identity of that genome inferred in the
 911 POInT computation. Losses from the triplicated state are then increasingly disfavored first to **D**_{1,3}
 912 (parameter $f_{1,3}$) and **D**_{2,3} (parameter $f_{2,3}$). There are also individual rates of loss from the duplicated to
 913 single-copy states (σ_x). In the null model, $f_{1,3}=f_{2,3}=1.0$ and $\sigma_1=\sigma_2=\sigma_3$. **WGQ**: Models of octoploid
 914 formation. The null model simply treats the four subgenomes as equivalent and as starting in the
 915 quadruplicated state (**Q**). This model has different loss rates from triplicated to duplicated loci (**T**_{x,y,z} to
 916 **D**_{x,y}, parameter δ) and duplicated to single-copy loci (**D**_{x,y} to **S**_x, parameter σ). A formation model for
 917 the octoploidy can then be added: all pillars start in state **D**_{1,3} and can symmetrically experience a
 918 gene loss from genome 1 or 3 (parameter λ) and transition to state **D**_{1,2} or **D**_{3,4} or become
 919 quadruplicated (null transition). The three models illustrated here are the most complex model fit to
 920 the various events, including the parameters associated and their numerical ranges.

921
 922 **Figure 3: A)** The assumed or computed phylogenetic relationships among species sharing the ten
 923 polyploidies studied (see *Methods*). Grey branches are those where no polyploidy event was studied.
 924 Because the temporal divergences of various groups are not well established, the tree is illustrated in
 925 an ultrametric format (Scaled topologies are shown in Supplemental Figure 1). Each polyploid branch
 926 is colored using POInT's estimates of the proportion of loci that were single-copy at its beginning and
 927 ending. Corresponding color keys for WGD, WGT and WGQ events are shown. The number of
 928 "pillars" (homoeologous loci) and the total number of gene models studied across each event are
 929 noted, as are the total number of loci and genes considered. The "*" on the yeast WGD branch
 930 indicates the branch where the proportion of genes returned to single-copy that are presently
 931 essential was tested (Supplemental Table 8). Next to each event, we show arrows and parameter
 932 estimates indicating post-polyploidy evolutionary processes such as biased fractionation for which we
 933 found significant evidence in that event (see key). **B)** An example mirrored gene tree for a completely
 934 retained set of homoeologs from At- α , illustrating the trees from which synonymous divergences were
 935 estimated. The branch lengths are given in number of synonymous substitution per synonymous site
 936 (i.e., K_s), with the shared internal (i.e., "root") branch shown in purple (K_s^R). For analysis purposes, the
 937 length of this branch was always divided by two to be comparable to the remaining branches (i.e.,
 938 split at its midpoint).

939
 940 **Figure 4:** Rapid loss of homoeologs immediately after polyploidy. On the x-axis is the ratio of
 941 rate of homoeolog loss (the α branch length estimate from POInT's models, see Figure 3) and
 942 the estimated mean synonymous divergence for that branch (\bar{K}_s ; see *Methods*). Hence, larger
 943 values of this ratio indicate more homoeolog losses per unit K_s . For the At- α , Brassica WGT,
 944 Legume WGD, Paramecium WGD and the TGD, the $\alpha \cdot t / \bar{K}_s$ ratio for the root branch is
 945 significantly larger than seen on any other branch (c.f., the 95% confidence intervals shown,
 946 computed as described in the *Methods* section). For these panels, we used a model excluding
 947 duplicate fixation here because including fixation in the model occasionally results in very long
 948 estimates of tip branch lengths (*Methods*). However, our conclusions are similar under the full
 949 WGD_{bfc-nb} model (see Supplemental Figure 3). On the y-axis is the net synonymous divergence to
 950 the end of the branch in question: in other words, the sum of the synonymous divergence of that
 951 branch and all its ancestors back to the root branch. This net divergence value is a rough
 952 indicator of the time since the polyploidy event for each branch. The root branch is indicated with
 953 a circle, other internal branches with squares and tip branches with triangles.

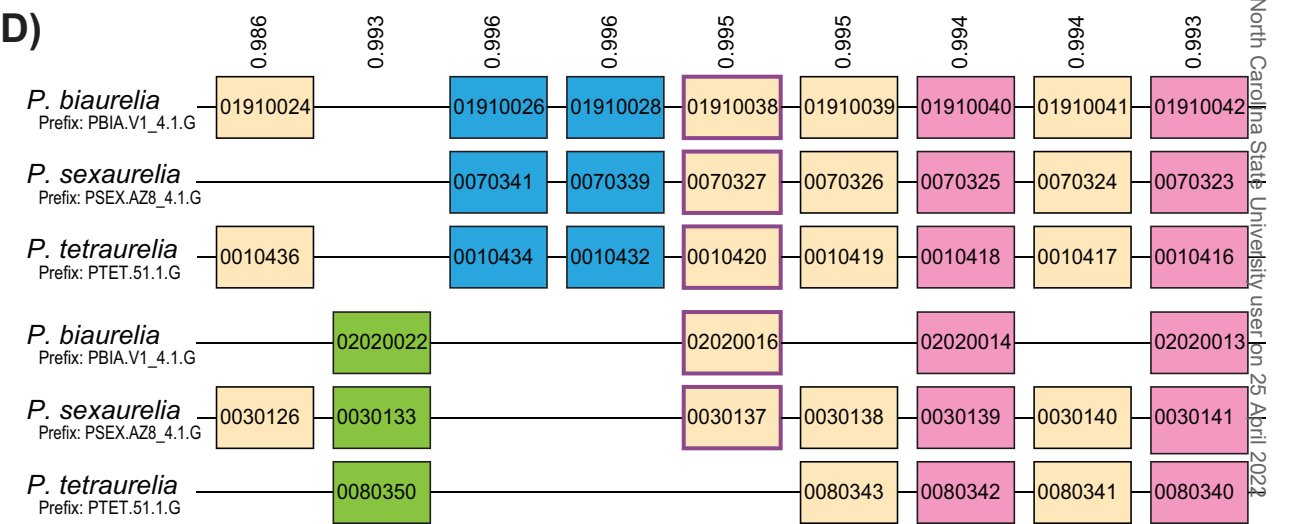
954
 955 **Figure 5:** Reciprocal gene loss (RGL) after polyploidy. **A)** Reciprocal gene losses (RGLs) between
 956 pairs of polyploid taxa (x-axis, normalized by the total number of loci/pillars analyzed for that event)

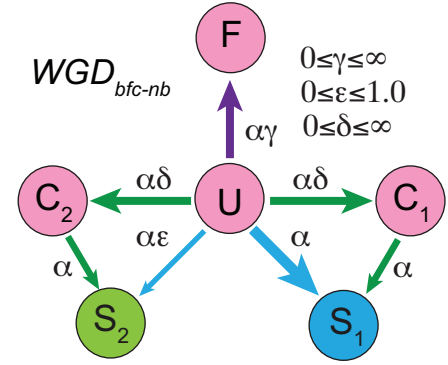
957 as a function of the inferred synonymous divergence of those taxa (*y*-axis). Panel **A** gives a
958 cropped view that focuses on RGLs in the non-yeast taxa, while panel **B** shows how the RGL
959 frequencies in yeast dramatically exceed those for the remaining events. For each pair of taxa from
960 a given event, we identified all single-copy loci in the two genomes where POInT infers a 95% or
961 greater confidence that those genes are paralogs created by the ancient polyploidy and not more
962 recent orthologs produced by the post-polyploidy speciation events. There are roughly linear
963 relationships between RGL frequency and synonymous divergence. Because the data points
964 shown are phylogenetically dependent (different species pairs share considerable common
965 evolutionary history), we have not attempted to fit regression lines to these data. Standard
966 approaches to phylogenetically independent contrasts (FELSENSTEIN 1985) do not apply here as the
967 inferred RGLs are pairwise species traits and not independent measures on each taxon. It is
968 however notable that the asexually reproducing yeasts appear to accumulate more RGLs per unit
969 K_s than other taxa. **B)** As for **A** but including the full range of RGL prevalence in the taxa sharing the
970 yeast WGD. **C)** Total numbers of RGLs inferred for each pair of taxa for each event (*x* axis).
971
972



C)

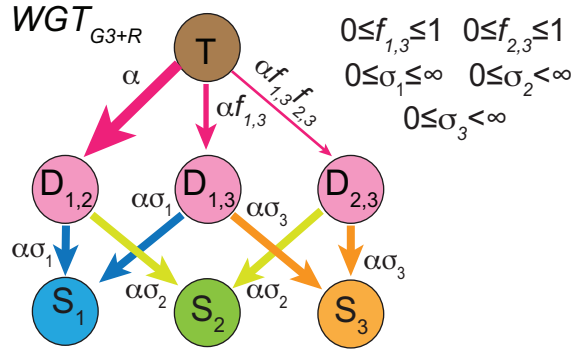
$$\underbrace{\begin{pmatrix} L_{00}^{i|D_i \dots D_0} \\ L_{01}^{i|D_i \dots D_0} \\ L_{10}^{i|D_i \dots D_0} \\ L_{11}^{i|D_i \dots D_0} \end{pmatrix}}_{L^{i|D_i \dots D_0}} = \underbrace{\begin{pmatrix} L_{00}^i \\ L_{01}^i \\ L_{10}^i \\ L_{11}^i \end{pmatrix}}_{L^i} \odot \underbrace{\begin{pmatrix} (1-\theta_i)^2 & \theta_i(1-\theta_i) & \theta_i(1-\theta_i) & \theta_i^2 \\ \theta_i(1-\theta_i) & (1-\theta_i)^2 & \theta_i^2 & \theta_i(1-\theta_i) \\ \theta_i(1-\theta_i) & \theta_i^2 & (1-\theta_i)^2 & \theta_i(1-\theta_i) \\ \theta_i^2 & \theta_i(1-\theta_i) & \theta_i(1-\theta_i) & (1-\theta_i)^2 \end{pmatrix}}_{\Theta} \cdot \underbrace{\begin{pmatrix} L_{00}^{i-1|D_{i-1} \dots D_0} \\ L_{01}^{i-1|D_{i-1} \dots D_0} \\ L_{10}^{i-1|D_{i-1} \dots D_0} \\ L_{11}^{i-1|D_{i-1} \dots D_0} \end{pmatrix}}_{L^{i-1|D_{i-1} \dots D_0}}$$





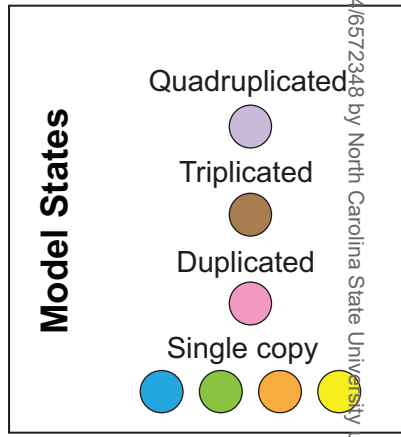
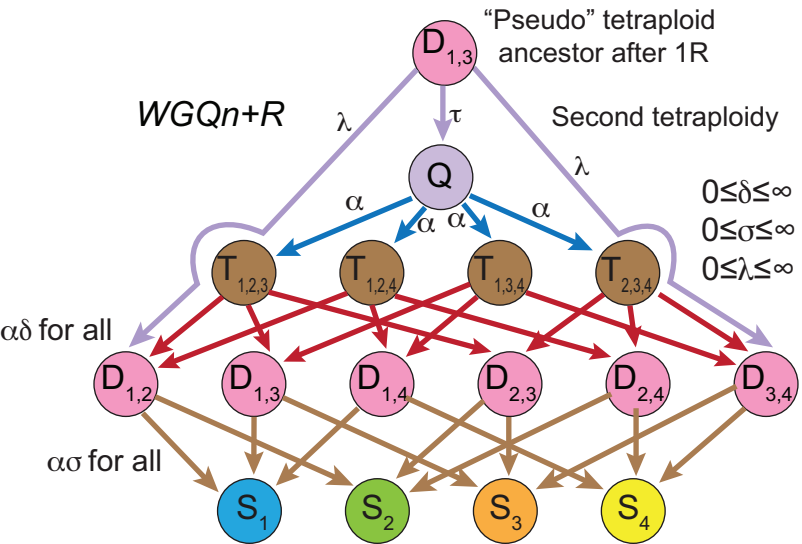
Convergent losses $\delta > 0$ \rightarrow Homoeolog fixation $\gamma > 0$

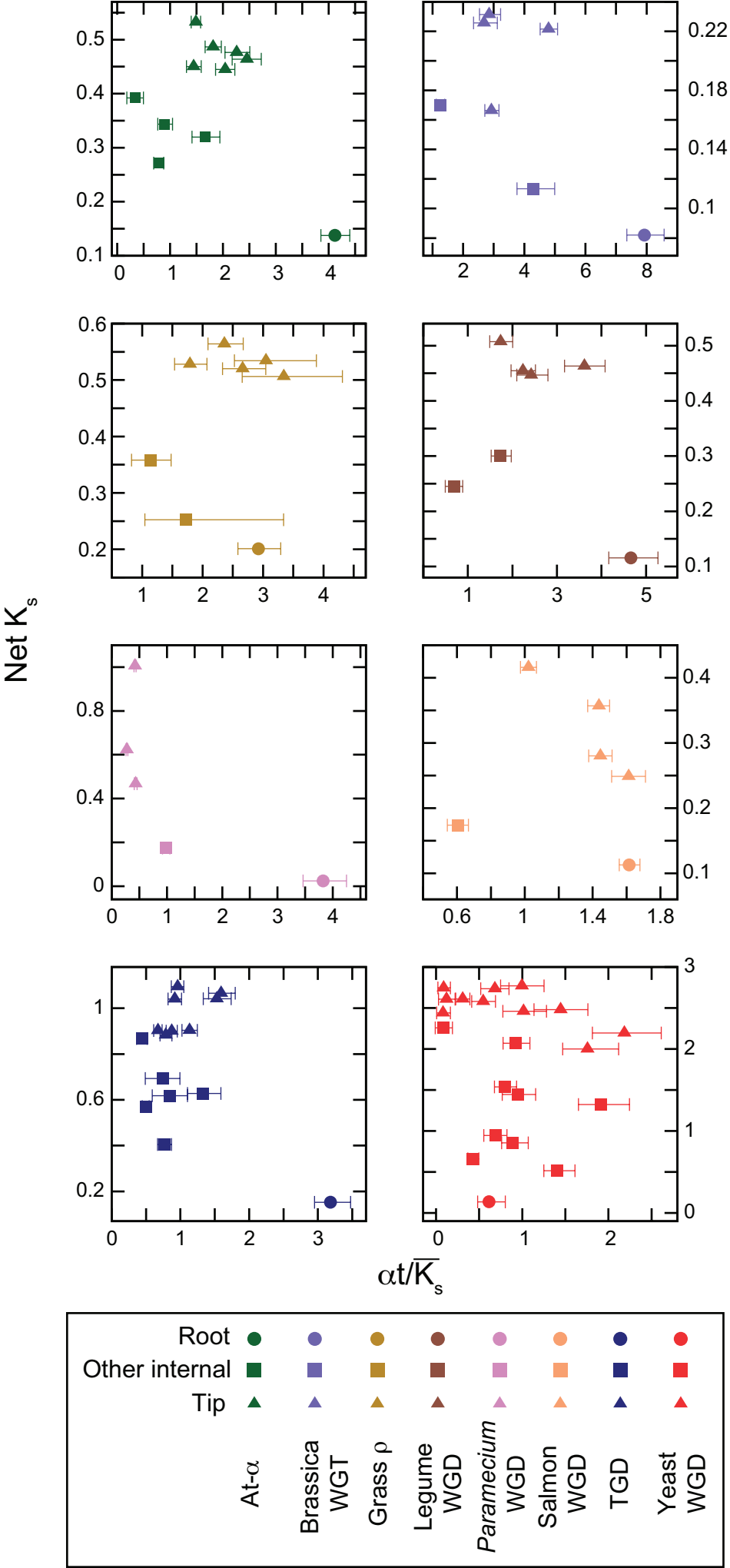
Biased fractionation $\epsilon \neq 1.0$



Biased fractionation $\rightarrow f_{1,3} \neq 1, f_{2,3} \neq 1$

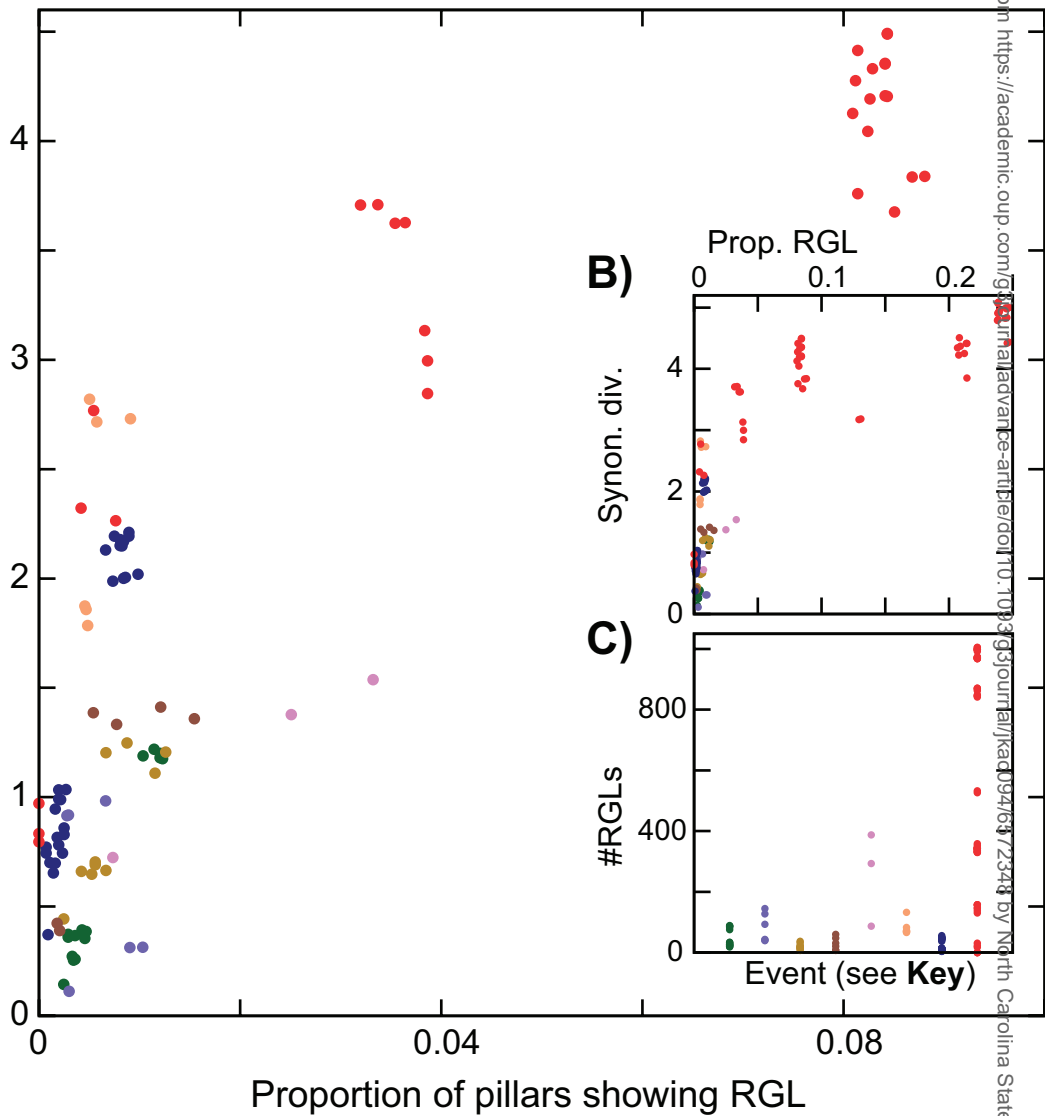
$\rightarrow \sigma_1 \neq \sigma_2$
 $\rightarrow \neq \sigma_3$



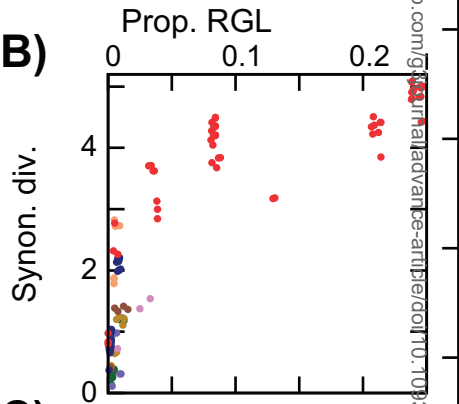


A)

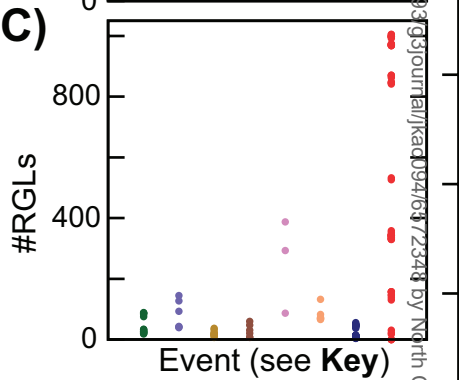
Synonymous divergence



B)



C)



Key

