

1 **Title: The genomic view of diversification**

2 Julie Marin¹, Guillaume Achaz^{1,2}, Anton Crombach^{1,3,4}, Amaury Lambert^{1,5}

3 ¹ *Center for Interdisciplinary Research in Biology (CIRB), Collège de France, CNRS UMR 7241,*
4 *INSERM UMR 1050, PSL Research University, Paris, France;*

5 ² *Institut de Systématique, Évolution, Biodiversité (ISYEB), MNHN, CNRS, Sorbonne Université,*
6 *EPHE, Paris, France;*

7 ³ *Inria, Lyon Antenne La Doua, Villeurbanne, France*

8 ⁴ *Université de Lyon, INSA-Lyon, LIRIS, UMR 5205, Villeurbanne, France*

9 ⁵ *Laboratoire de Probabilités, Statistique et Modélisation (LPSM), Sorbonne Université, CNRS UMR*
10 *8001, Paris, France.*

11

12 **ABSTRACT:** Evolutionary relationships between species are traditionally represented in the form of a
13 tree, called the species tree. The reconstruction of the species tree from molecular data is hindered
14 by frequent conflicts between gene genealogies. A standard way of dealing with this issue is to pos-
15 tulate the existence of a unique species tree where disagreements between gene trees are explained
16 by incomplete lineage sorting (ILS) due to random coalescences of gene lineages inside the edges of
17 the species tree. This paradigm, known as the multi-species coalescent (MSC), is constantly violated
18 by the ubiquitous presence of gene flow revealed by empirical studies, leading to topological incon-
19 gruences of gene trees that cannot be explained by ILS alone. Here we argue that this paradigm
20 should be revised in favor of a vision acknowledging the importance of gene flow and where gene
21 histories shape the species tree rather than the opposite. We propose a new, plastic framework for
22 modeling the joint evolution of gene and species lineages relaxing the hierarchy between the species
23 tree and gene trees. [We implement this framework in two mathematical models called the gene-based](#)
24 [diversification models \(GBD\): 1\) GBD-forward, following all evolving genomes and thus very intensive](#)
25 [computationally and 2\) GBD-backward, based on coalescent theory and thus more efficient. Each](#)
26 [model features four parameters tuning colonization, mutation, gene flow and reproductive isolation.](#)
27 [We propose a quick inference method based on the differences between gene trees and use it to eval-](#)
28 [uate the amount of gene flow in two empirical data-sets. We find that in these data-sets, gene tree](#)
29 [distributions are better explained by the best fitting GBD model than by the best fitting MSC model.](#)
30 This work should pave the way for approaches of diversification using the richer signal contained in
31 genomic evolutionary histories rather than in the mere species tree.

32

33 **Keywords:** coalescent theory, gene flow, gene tree, gene-based diversification model, multi-species
34 coalescent, phylogeny, population genetics, speciation, species tree, reproductive isolation, introgres-
35 sion.

36

37 The most widely used way of representing evolutionary relationships between contemporary species
38 is the so-called species tree, or phylogeny. The high efficiency of statistical methods using sequence
39 data to reconstruct species trees, hence called ‘molecular phylogenies’, led to precise dating of the
40 nodes of these phylogenies [35, 38, 87]. Notwithstanding the debatable accuracy of these datings,
41 the use of time-calibrated phylogenies, sometimes called ‘timetrees’ [34], has progressively overtaken
42 a view where phylogenies merely represent tree-like relationships between species in favor of a view
43 where the timetree is the exact reflection of the diversification process [61, 70, 85]. In this view, the
44 nodes of the phylogeny are consequently seen as punctual speciation events where one daughter
45 species is instantaneously ‘born’ from a mother species. In this paper, we explore an alternative
46 view of diversification, acknowledging that speciation is a long-term process [17, 43, 72] and not
47 invoking any notion of mother-daughter relationship between species as done in the timetree view.
48 This alternative view is gene-based rather than species-based, comparable with Wu’s genic view of
49 speciation [91]. We use here the term ‘gene’ in the sense of “non-recombining locus”, *i.e.*, a region of
50 the genome with a unique evolutionary history. Our view is meant in particular to accommodate the
51 well-recognized existence of gene flow between incipient species, which persists during the speciation
52 process and long after [51].

53 The timetree view of phylogenies does acknowledge that gene trees are not independent and may
54 disagree with the species tree [48], but current methods jointly inferring gene trees and species tree
55 rely on the following assumptions that we question in the next section: there is a unique species
56 tree, the species tree shapes the gene trees and the species tree is the only factor mediating all
57 dependences between gene trees (they are independent conditional on the species tree).

58 This view is materialized in a model called the ‘multispecies coalescent’ (MSC) [39] where con-
59 ditional on the species tree, the evolutionary histories of genes follow independent coalescents con-
60 strained to take place within the hollow edges of the species tree. Many methods have been devel-
61 oped to estimate the species tree under the MSC, such as full likelihood methods (e.g. BEAST [35],
62 BPP [94]) which average over gene trees and parameters [93], and the approximate or summary co-
63 alescent methods (e.g. ASTRAL [58], MP-EST [45], and STELLS [92]) which use a two-step approach:
64 gene trees are first inferred and then combined to estimate the species tree that minimize conflicts
65 among gene trees. Discordance between gene topologies is then explained, as a first approximation
66 at least, by the intrinsic randomness of coalescences resulting in incomplete lineage sorting (ILS)
67 (figure 1).

68 However, the presence of gene flow (introgression, hybridization, horizontal transfer) is now widely
69 recognized between closely related species, and even between distantly related species [51]. Porous
70 species boundaries, allowing for gene exchange because of incomplete reproductive isolation, are
71 indeed regularly observed in diverse taxa such as amphibians [21, 68], arthropods [12], cichlids [90],
72 cyprinids [6, 24, 25, 26, 84], insects [63, 67, 89], and even more frequently among bacteria [51, 83].

73 Long neglected, gene flow has recently been recognized as an important evolutionary driving force,
74 through adaptive introgression or the formation of new hybrid taxa [1]. The ubiquity of genetic ex-
75 change across the Tree of Life between contemporary species suggests that gene flow has occurred
76 many times in the evolutionary past, and might actually be the most important cause of discrepan-
77 cies between gene histories (e.g. [8, 11, 23, 37]) (figure 1). Accordingly, several extensions to the
78 MSC model have been considered allowing for gene flow between species [40, 95]. These models ac-
79 knowledge that species boundaries can be permeable at a few specific timepoints [33]. Unfortunately,
80 because of the heavy computational cost of modeling the coalescent with gene flow, these methods
81 are limited to small data-sets [95]. More importantly, they might not be appropriate to realistically
82 model gene flow, given the frequency of gene flow across time and clades described in empirical
83 studies [82]. Additionally, some of these methods, ASTRAL and MP-EST, might infer erroneous gene
84 trees when gene flow is present [47]. These observations urge for novel approaches where gene flow
85 is the rule rather than the exception.

86 To fill this void, we propose here [an alternative framework and two accompanying models \(one](#)
87 [in forward time and one in backward time\)](#), the gene-based diversification (GBD) models, framed
88 with minimal assumptions arising from recent empirical evidence. [Those models rely on the property](#)
89 [of populations to spontaneously differentiate genetically \(mutation\) while simultaneously undergoing](#)
90 [gene flow. This genetic differentiation is accompanied by a decrease in gene flow until reproductive](#)
91 [isolation is complete \(these processes are detailed below\).](#) Moreover, unlike previous models, we
92 [place ourselves in the case of pervasive gene flow among species that may have occurred countless](#)
93 [times in the past, as suggested by recent studies. The GBD models are anchored in a new concep-](#)
94 [tual framework, that we call the genomic view of diversification.](#) Unlike the timetree view, the present
95 framework does not put the emphasis on the species tree (which in our model becomes a network
96 rather than a tree) and assumes that gene trees shape the species tree (rather than the opposite).

97

98 **THE GENOMIC VIEW OF DIVERSIFICATION**

99 **Gene flow and the questionable existence of a species genealogy**

100 The biological species concept (BSC [54]) defines species as groups of interbreeding populations
101 that are reproductively isolated from other groups. This definition postulates the non-permeability of
102 species boundaries, which is contradicted by the growing body of evidence describing permeable or
103 semi-permeable genomes, even between distantly related taxa. To integrate the possibility of gene
104 flow into the definition of species, Wu [91] shifted the emphasis from isolation at the level of the whole
105 genome to differential isolation at the gene level. Species are thus defined as differentially adapted
106 groups for which inter-specific gene flow is allowed except for genes involved in differential adaptation
107 (a well-defined form of divergence in which the alternative alleles have opposite fitness effects in
108 the two groups) [91]. Because a fraction of the genome may still be exchanged after speciation

109 is complete, a mosaic of gene genealogies is expected between divergent genomes [91]. Much
110 evidence supports this prediction with the observation of highly conflicting gene trees, e.g. Darwin's
111 finches [27, 29], sympatric sticklebacks [73, 77], Iberian barbels [25], and *Rhagoletis* species [2].

112 Accordingly, the notion of a species genealogy as the binary division of species into new inde-
113 pendently evolving lineages in bifurcating phylogenetic trees, appears inappropriate. To avoid this
114 misleading vision of speciation, we here wish to relax the species tree constraint by considering
115 only gene genealogies as real genealogies, thereby laying aside, at least temporarily, the notion of
116 species genealogy. To do so, we do not specify mother-daughter relationships between species, yet
117 we postulate the existence of species at any time, and assume that we can unambiguously follow the
118 genealogies of genes (defined as non-recombining loci, as mentioned above).

119 Looking forward in time, genes belonging to two distinct individuals may find each other, in a next
120 generation, in the same genome because of recombination. The same process might occur with two
121 individuals belonging to different species under gene flow.

122 The notion of a species genealogy as a binary bifurcating tree is hardly compatible with gene
123 flow, and a direct consequence is to challenge the notion of a unique ancestral species. If all genes
124 ancestral to species S have travelled through the same species in the past, then species S has only
125 one single ancestor species at any time. But because of gene flow, these genes may lie in different
126 species living at a given time in the past, such that species S can have several ancestral species at
127 this time. In other words, several species have contributed to the present-day genome of the species
128 S .

129

130 **Genomic coadaptation under continuous gene flow**

131 While some genes (e.g., genes involved in divergent adaptation) are hardly exchanged between
132 populations, other genes (e.g., neutral genes unlinked to genes under divergent selection) can be sub-
133 ject to gene flow between different species [69, 91]. Gene flow can persist for long periods of time, with
134 evidence suggesting introgression events occurring over periods lasting up to 20 Myr [6, 25, 90]. Over
135 time, genetic differences will accumulate in regions of low recombination and expand via selective
136 sweeps, leading eventually to complete reproductive isolation [91]. Because populations differentially
137 accumulate new alleles, their compatibility (hybrid fitness) will be affected. This process has been con-
138 ceptualized by Dobzhansky and Muller [13, 62] in the so-called Bateson-Dobzhansky-Muller (BDM)
139 model [10]. This model proposes that genetic incompatibilities, hence called BDM incompatibilities,
140 are characterised by negative epistatic interactions between alleles at two or more genes that have
141 fixed differentially, in each of the parental populations, by local adaptation or genetic drift. The selec-
142 tive value of hybrids is reduced because the new alleles, divergently selected in each populations, are
143 not adapted to each other. On the other hand, in the parental populations these alleles are co-adapted
144 and have neutral or even beneficial effects [79, 88]. These incompatibilities have been hypothesized

145 to increase at a rate proportional to the square of time [64]. Accordingly, pairs of species will likely
146 exhibit greater genetic incompatibility as a function of time since divergence, *i.e.* be less permeable
147 to gene flow, as has been observed for Iberian barbels [25], pea aphids [67], or salamanders [68]. In
148 other words, gene lineages remaining too long isolated within different species decrease their ability
149 to introgress the genome of the other, a property that we name *genomic coadaptation* and which is
150 the consequence of spontaneous mutation.

151

152 **The gene-based diversification (GBD) models**

153 We propose here a new plastic framework, derived from the genomic view of diversification de-
154 scribed above, that acknowledges the importance of gene flow and relaxes the hierarchy between the
155 species tree and gene trees. We built two models, one in forward time that follows the standard view
156 of the main biological processes responsible for diversification under gene flow, and one in backward
157 time, less computationally intensive, with matching backward parameters (figure 2). These models
158 that we named the gene-based diversification (GBD-forward and GBD-backward) models, use coales-
159 cent theory for modelling the joint evolution of gene and species lineages, reconciling phylogenomics
160 with our current knowledge of species diversification. The biological mechanisms first, then the cor-
161 responding parameters, are detailed thereafter for each model.

162

163 *The GBD-forward model*

164 The GBD-forward model describes the joint action of four processes affecting the diversification of
165 genomes (see figure 2): colonization, mutation, drift and introgression.

166 We consider a stochastically varying number of *populations*, all populated with individual genomes.
167 We neglect extinctions and focus on *colonization* events, at which one population seeds a daughter
168 population founded by one or several of its individuals. Genes independently accumulate *mutations*
169 with time, under the infinite-allele model assumption. Mutations can be fixed or lost due to selection
170 and genetic drift, that we summarize here under the term *drift*.

171 As a result of mutations and drift, populations differentiate genetically through time, which results
172 in the decrease of gene flow. To model this, we follow what we term the *co-adaptation* between non-
173 homologous genes and assume that *introgression* is governed by the numbers of co-adapted alleles
174 in the receiver and donor populations. Right after colonization, all the genes of the daughter and
175 mother populations carry the same alleles and so are co-adapted. Now an allele having arisen at
176 time t by mutation on some gene is co-adapted only with the alleles carried by its genome at time t .
177 This assumption underlies the well-known model of BDM incompatibilities described previously. Each
178 time a mutation occurs the number of co-adapted genes among populations will decrease, reducing
179 in turn the possibility of genetic exchange between populations.

180 Two populations that are completely differentiated, in the sense that all pairs of non-homologous

181 alleles sampled from each of them are not co-adapted, can no longer exchange genes and can thus
182 be seen as different species. Because populations are constantly differentiating from each other, we
183 name populations in the prospective point of view (GBD-forward) what will become species only from
184 a retrospective point of view (GBD-backward).

185 Demographic events are assumed to be much faster than other processes. In the time scale con-
186 sidered here, (1) the fixation of alleles within populations is instantaneous so that all genomes in a
187 population are identical (we thus do not model the co-existence of several different homologous alle-
188 les within a population) and (2) a colonization event can be seen as the instantaneous replication of
189 one population into two, actually because of (1), of one genome into two.

190

191 Parametrization

192 At $t = 0$, we consider a single monomorphic population, summarized into a single genome har-
193 boring n genes. During the diversification process, the genome of this population (n genes) will be
194 replicated, mutations will be differentially fixed in each population, and the genomes of these popula-
195 tions can be replicated again. We follow the lineages of these n genes in forward time, assuming a
196 time-discrete Markov chain associated to the time-continuous chain with the following rates.

- 197 • **Mutation** (rate α). At any time t , each gene lineage in each population can acquire a new allele
198 (infinite-allele model) at rate α . By definition, a new allele occurring at gene L on genome G
199 is co-adapted with the allele present at a gene L' , for any L' (different of L) of genome G . On
200 the contrary, a mutation arising at gene L of genome G and a mutation arising at gene L' of
201 genome G' are not co-adapted.
- 202 • **Colonization** (rate β). At any time t , each population can be replicated at rate β into a new
203 population which will evolve independently in the future. The newborn population is assumed to
204 carry the same genome as carried by the mother population.
- 205 • **Genetic drift** (rate γ). Each population undergoes Moran-type births and deaths at rate γ . In
206 this work, we assume γ to be much larger than all other parameters, so that each population is
207 actually monomorphic at all times.
- 208 • **Introgression** (rate δ). At any time t , each gene lineage at locus L on genome G can be repli-
209 cated and introgress genome G' at rate $\delta(n-1)$, proportional to the number of non-homologous
210 loci in genome G' . If accepted by the target genome G' , the replicated lineage replaces its ho-
211 mologous gene lineage (at locus L in G'). The introgression is accepted with a probability equal
212 to the fraction of the $n-1$ non-homologous genes on G' carrying an allele co-adapted with the
213 allele carried by L .

214 Diversification occurs until a number K of different populations is reached and the whole process is

215 stopped when the K populations are genetically isolated, that is, when no pair of alleles carried by
216 different genomes is co-adapted (i.e., when all probabilities of introgression are equal to 0).

217 This framework can be made more complex by letting the parameters depend on time, on the
218 gene, or on any prescribed category of genes.

219

220 *The GBD-backward model*

221 The GBD-backward model is not the exact backward picture of the GBD-forward model but relies
222 on the same idea that genomes in different populations tend to diverge with time until they cannot
223 exchange alleles. The consequence of this fact is that genes sampled in the same genome today will
224 tend to be found in the same population in the past more often than by chance. We model this phe-
225 nomenon by saying that the ancestral lineages of genes sampled in the same present-day genome
226 are *co-adapted*, and that co-adapted genes are *attracted* towards each other. The GBD-backward
227 model describes the joint action of four processes (see figure 2): non-homologous attraction, homol-
228 ogous attraction, coalescence and gene flow.

229 As explained above, in the retrospective point of view (GBD-backward), we name species the
230 populations in which the ancestral gene lineages travel.

231 When two homologous gene lineages are in the same species they can *coalesce* when finding
232 their common ancestor, that is merge into a single lineage (hence within the same genome).

233 Each gene lineage can move from its species to another species. This happens as a result of
234 homologous attraction, non-homologous attraction and gene flow. As explained previously, (non-
235 homologous) *co-adapted* genes move into the same species as a result of *non-homologous attraction*,
236 which can be viewed as the backward consequence of *mutations*. Homologous gene lineages move
237 into the same species as a result of *homologous attraction*, which can be viewed as the backward
238 picture of a *colonization event*, when populations and their genomes have been replicated. Last, any
239 gene lineage can move from its species by *gene flow* to an empty species, i.e., a species containing
240 no other gene lineage ancestral to the sample.

241 Note that after coalescence of two homologous lineages, the resulting lineage is now ancestral to
242 at least two genomes and thus co-adapted with all gene lineages ancestral to these genomes. As
243 a consequence of the mere *non-homologous attraction*, going further back in time, all other genes
244 will then move to the same species and further coalesce, until all homologous gene lineages have
245 coalesced.

246 Equivalently to the *drift* process in forward time, we will assume that the *coalescences* are fast, so
247 that in backward time homologous attraction events are immediately followed by coalescence of the
248 two gene lineages.

249

250 Parameterization

251 At $t = 0$, n homologous genes are sampled in each of N distinct species. Retrospectively, the
252 genomes of these N species (harbouring each n genes) will merge progressively in one genome
253 of n genes at some time t in the past. Homologous genes, one by one, will merge (*homologous*
254 *attraction and coalescence*). Merged genes will then attract all the genes of their original genomes
255 (*non-homologous attraction*), until the *coalescence* of all homologous genes. We follow the lineages
256 of these n genes in backward time, assuming a time-discrete Markov chain associated to the time-
257 continuous chain with the following rates.

- 258 • **Non-homologous attraction** (rate a). At any time t in the past, as a backward picture of
259 **genomic coadaptation**, each gene lineage L escapes from its species S at rate $a(n - 1)$ per
260 target species S' , proportional to the number of non-homologous loci in the genome G' hosted
261 by S' . It is accepted in S' based on its co-adaptation with G' . Recall that if G_0 denotes the
262 genome harboring the descendant lineage of L at time $t = 0$, then all gene lineages harbored
263 by G' that are ancestral to G_0 are said co-adapted with L . Then L is accepted in S' with a
264 probability proportional to the fraction of the $n - 1$ non-homologous loci of G' that are co-adapted
265 with it. The parameter a corresponds to the *mutation* parameter α of the GBD-forward model.
- 266 • **Homologous attraction** (rate b). At any time t in the past, each gene lineage at rate b per
267 homologous gene lineage, moves to the species harboring this homologous lineage (or in an
268 alternative, more specific version of the model, each gene lineage belonging to some previ-
269 ously prescribed category, like genes contributing to reproductive isolation). This parameter
270 corresponds to the *diversification* parameter β of the GBD-forward model.
- 271 • **Coalescence** (rate c). At any time t in the past, each pair of homologous genes lying within the
272 same species coalesces at rate c . This parameter corresponds to the *genetic drift* parameter γ
273 of the GBD-forward model.
- 274 • **Gene flow** (rate d). At any time t in the past, as a backward picture of introgression, each
275 gene lineage escapes from its genome at rate d and enters an empty species (also called
276 ghost species, *i.e.*, harboring no other gene lineage ancestral to the samples, figure 2). This
277 parameter corresponds to the *introgression* parameter δ of the GBD-forward model. To model
278 the introgression of bigger chunks of DNA, we could alternatively assume that instead of one
279 lineage, a given fraction of the lineages of a genome can simultaneously move to an otherwise
280 empty species. We will not consider this possibility in the present work.

281 We define the number of ancestral species of a given genome at time t , as the number of species at
282 time t containing gene lineages ancestral to this genome. We considered a *time unit* to be equal to
283 the time elapsed between two events that we assumed to be constant for the sake of simplicity. In this
284 manuscript we wish to explore the impact of gene flow rather than ILS to explain gene tree conflicts,

285 and thus consider a large c value (coalescence rate) so that coalescence events are instantaneous,
286 which is consistent with the large γ value of the forward model. Therefore, only the parameters a , b ,
287 and d have an influence on the gene genealogies in the GBD-backward model.

288

289 The GBD models were implemented in R (<https://www.r-project.org>) and evaluated under different
290 sets of parameters. Because the GBD-forward model is computationally prohibitive, while giving com-
291 parable results with the GBD-backward model, we conducted most of the analyses and the inferences
292 with the GBD-backward model. We provide a preliminary, ABC-like inference method by minimizing
293 the difference (Kullback-Leibler divergence) between the distributions of Billera-Holmes-Vogtmann
294 (BHV) distances (pairwise distances between gene trees) [3] in empirical vs simulated data. We
295 applied this inference method to two empirical multi-locus data-sets showing complex evolutionary
296 patterns due to gene flow, each comprising six morphologically and ecologically distinct species, the
297 Ursinae (a bear subfamily) [42] and the *Geospiza* clade (a genus of Darwin's finches) [20]. We es-
298 timated in particular 1) the relative amount of gene flow that has shaped each data-set, and 2) the
299 corresponding average number of ancestral species.

300

301 MATERIAL AND METHODS

302 Inference method for the GBD-models

303 When considering several sampled genomes all containing n genes, a set of n gene trees is
304 obtained for each particular parameter setting and each realization of the model. To characterize a
305 set of gene trees, we employed a multidimensional summary statistic defined as the distribution of
306 pairwise distances between gene trees. Because the GBD-models are time oriented, a tree metric for
307 rooted trees was necessary. Among this class of metrics, those accounting only for topology (such
308 as the Robinson–Foulds metric [71]) reached a plateau for large amounts of gene flow because the
309 maximal distance among all pairs of gene trees was reached (results not shown). For these reasons
310 we opted for the Billera-Holmes-Vogtmann (BHV) metric [3] that accounts for both branch lengths and
311 topological differences, allowing us to distinguish sets of gene trees even when simulated with much
312 gene flow. This metric is based on a view of tree space as a quadrant complex with quadrants sharing
313 faces. Two trees with the same topology lie in the same quadrant, otherwise they lie in two distinct
314 quadrants. At a common edge between two quadrants, the incongruent internal branches between
315 trees have lengths equal to zero. Then a distance can be calculated between two rooted trees as the
316 shortest path across these interconnected quadrants.

317 BHV distances do not rely only on the topology but also on branch lengths. The difference in topol-
318 ogy is weighted by the branch lengths supporting these topologies, therefore uncertainties causing
319 polytomies (or a branching pattern close to a polytomy) in gene trees will only marginally affect our
320 results.

321 To compare trees that did not evolve on the same time scale, BHV distances were computed on
322 re-scaled trees. For each set of gene trees issued from a single simulation or data-set, we rescaled
323 all the trees so that the median of the most recent node depth is 1.

324 To find the best set of parameters $(\frac{1}{a}, b)$, for empirical or test (simulated) trees, we employed the
325 Kullback-Leibler (KL) divergence (package 'FNN' in R) as a distance metric by minimizing this dis-
326 tance between the distributions of BHV pairwise distances of empirical, and test trees, with simulated
327 trees. The lower the KL divergence is the better is the fit.

328

329 **Inference method accuracy**

330 We tested the accuracy of our inference method, minimization of the difference (KL divergence)
331 between the distributions of BHV distances, on 8 simulated data sets (test trees). Using the GBD-
332 backward model, we built gene trees for 8 sampled combination parameters $(\frac{1}{a}, b)$ with $\frac{1}{a} \in [0.3, 3.5]$,
333 every 0.2, and $b \in [0.01, 0.12]$, every 0.01. The other parameters were fixed, $d = 1$, $c = 200$, $n = 10$
334 and $N = 6$. The number of time units t was set to 5,000, which guarantees the coalescence of all
335 homologous genes. We performed 15 replicates.

336 We next optimized the GBD-backward model for $N = 6$ and $n = 10$ by varying two parameters,
337 a and b , and fixing $d = 1$ and $c = 200$. The number of time units t was set to 5,000. We performed
338 15 replicates under each parameter combination in a grid of $(\frac{1}{a}, b)$ with $\frac{1}{a} \in [0.3, 3.5]$, every 0.2, and
339 $b \in [0.01, 0.12]$, every 0.01. The parameters $(\frac{1}{a}, b)$ of the test trees were inferred by minimizing the KL
340 distance between the test trees and the simulated trees from the grid.

341

342 **Comparison of the GBD-models**

343 To visually compare the reconstructed genealogies obtained with the GBD-forward and the GBD-
344 backward model we performed simulations for genomes containing $n = 5$ genes, with $\alpha = 0.5$, $\beta = 1$,
345 $\delta = 0.2$ and $K = 30$ for the GBD-forward, and $a = 1$, $b = 0.1$, $d = 2$ and $N = 10$ for GBD-backward
346 model.

347 Next, we used our inference method (minimization of the KL distances between distributions of
348 BHV pairwise distances) to compare the parameters of the two GBD models. We simulated the GBD-
349 forward model with the following parameters: $\alpha = 0.5$, $\beta = 0.01, 0.02$ and 0.05 , $\delta \in [0.01, 0.04]$, every
350 0.01 , $n = 10$ and $K = 30$. For each set of parameters, 6 replicates were performed and averaged.
351 The simulations were stopped when the number of populations reached $K = 30$, and the trees were
352 built from the first $N = 6$ genomes (populations). The KL distances were then minimized between
353 the distributions of BHV pairwise distances of the GBD-forward trees and the inference grid obtained
354 from the GBD-backward model and described in the previous section.

355 In order to compare the GBD-forward and GBD-backward trees, we only took into account the
356 colonization and introgression events affecting the $N = 6$ genomes when reconstructing the GBD-

357 forward trees. Indeed in forward time, all the genomes (here $K = 30$) are simulated with all the
358 corresponding events affecting their genes (mutation, colonization and introgression). In backward
359 time the events for only $N = 6$ genomes are simulated. Moreover in backward time mutations are not
360 modeled, whereas in forward time each new mutation corresponds to a distinct event.

361 Both models gave qualitatively similar results (see the results section). However because the
362 GBD-forward model was computationally prohibitive, all the following analyses were performed with
363 the GBD-backward model. A simulation, with $N = 6$ ($K = 30$ for the GBD-forward to be able to recon-
364 struct genealogies of $N = 6$ genomes), $n = 10$, $a/\alpha = 1$, $b/\beta = 1$, $c = 200$ and $d/\delta = 1$, took about
365 10 hours for the GBD-forward model and 10 minutes for GBD-backward model (Intel(R) Core(TM)
366 i7-6700 CPU).

367

368 **A single sampled genome (GBD-backward model)**

369 We aimed to evaluate the variation in the number of ancestral species with gene flow. We per-
370 formed simulations for a single sampled genome containing n genes (with $n = 20, 50, 100, 200$),
371 and varied the relative amount of introgression (*gene flow*) compared to genetic differentiation (*non-*
372 *homologous attraction*), ratio $\frac{d}{a}$ (with $a = 1$ and $d \in [0.2, 2]$, every 0.2). The number of time units t
373 was set to 10,000. We sampled the number of ancestral species every 500 time units starting at time
374 $t = 5,000$, and averaged these 11 values for each simulation. For each set of parameters, 5 replicates
375 were performed and averaged.

376 A model is said to be *sampling consistent* if the same outcome is expected for any k sampled
377 genes independently of the total number n of genes in the genome. To evaluate the validity of this
378 property, we randomly sampled $k = 20$ genes from each genome of $n \geq 20$ genes and computed their
379 average number of ancestral species.

380

381 **A sample of several genomes (GBD-backward model)**

382 We evaluated the influence of the number n of genes (with $n = 5, 10, 20$), of the number of species
383 N (with $N = 6, 10$), and of the relative amount of gene flow $\frac{d}{a}$ (with $d = 1$ and $\frac{1}{a} = 0.3, 0.5, 0.9, 1.3, 1.7,$
384 $2.1, 2.5, 2.9, 3.3$) on gene tree diversity (BHV distances) (figure 6A). The other parameters were fixed,
385 with $b = 0.05$ and $c = 200$.

386 For the same values of $\frac{d}{a}$ and c , and for $n = 10$, $N = 6$, we also evaluated the influence of the
387 *homologous attraction* rate b (with $b = 0.01, 0.02, 0.05, 0.12$) on gene tree diversity (BHV distances)
388 (figure 6B).

389

390 **The GBD-backward model versus the MSC model**

391 To evaluate the ability of MSC methods to deal with gene flow, we estimated a species tree and its
392 gene trees (MSC model with no gene flow) using sequences corresponding to gene trees simulated

393 under the GBD-backward model (with gene flow).

394 A set of 10 gene trees was simulated under the GBD-backward model (with $N = 6$, $b = 0.05$, $\frac{1}{a} =$
395 0.9 , and $d = 1$) (figure 7). We simulated DNA sequences (package 'PhyloSim' in R [80]) corresponding
396 to each of the 10 gene trees with model of DNA evolution estimated by modeltest (function 'modelTest',
397 package 'phangorn' in R [76]) for the TRAPPC10 intron of the bear data-set detailed below [42]:
398 HKY model, rate matrix: $A = 1.00$, $B = 5.29$, $C = 1.00$, $D = 1.00$, $E = 5.29$, $F = 1.00$, base
399 frequencies: 0.26, 0.19, 0.21, 0.34. Prior to simulating the sequences, the 10 gene trees were scaled
400 to the TRAPP10 intron phylogenetic tree length (built with RaXML 8.1.11 [86] assuming GTR (general
401 time reversible) model with 1,000 bootstrap replicates).

402 The species tree and the gene trees associated were estimated from the simulated sequences
403 with the program BEAST v. 2.4.8 [5] with the following parameters: unlinked substitution models, un-
404 linked clock models, unlinked trees, HKY substitution model for each of the 10 genes, strict clock, Yule
405 process to model speciation events, and 80 million generations with sampling every 5000 generations.
406 To set the calibration time of the root we assumed that 1 time unit corresponded to 10 ky; on average
407 the last coalescence event among the 10 GBD-backward trees occurred at $t = 700$. Accordingly, we
408 used a normal distribution prior for the root heights (mean=7.0 (My); stdev=1.0).

409

410 Inference from empirical data-sets

411 *Empirical data-sets*

412 To evaluate if the GBD-backward model correctly reproduces the signal left by gene flow in gene
413 trees we simulated gene trees under the GBD-backward model and under the MSC model. The ad-
414 equacy between the simulated trees and empirical gene trees was estimated by comparing the dis-
415 tributions of pairwise gene tree distances of simulated vs empirical data-sets. The empirical clades
416 have been chosen for their moderate phylogenetic depth, good sampling coverage and known con-
417 flicting gene trees. The first data-set comprised 14 autosomal introns for 6 bear species (*Helarctos*
418 *malayanus*, *Melursus ursinus*, *Ursus americanus*, *U. arctos*, *U. maritimus*, and *U. thibetanus*) and 2
419 outgroups (*Ailuropoda melanoleuca* and *Tremarctos ornatus*) [42]. The sequences were downloaded
420 from GenBank (supplementary table S1). As in Kutschera et al. [42], all variation within and among
421 individuals was collapsed into one single 50% majority-rule-consensus sequence for each of the 8
422 species. The phylogenetic trees were built with the program BEAST v. 1.8.3. [14], with the param-
423 eters used by the authors of [42]: Yule prior to model the branching process, strict clock, a normal
424 prior on substitution rates (0.001 ± 0.001) (mean \pm SD), minimum age of 11.6 My for the divergence
425 of *A. melanoleuca* from other bears (exponential prior: mean= 0.5; offset= 11.6), and 10 million
426 generations with sampling every 1000 generations. The models of DNA evolution were estimated
427 by modeltest (function 'modelTest', package 'phangorn' in R [76]) (supplementary table S2). The
428 monophyly of the ingroup and the topology among the outgroups were constrained according to the

429 topology depicted in Kutschera et al. [42].

430 The second data-set comprised 7 nuclear markers for 6 finch species (*Geospiza conirostris*,
431 *G. fortis*, *G. fuliginosa*, *G. magnirostris*, *G. scandens*, and *G. septentrionalis*) and 2 outgroups (*Ca-*
432 *marhynchus psittacula* and *Platypiza crassirostris*) [20]. The sequences were downloaded from
433 GenBank (supplementary table S3). The phylogenetic trees were built with the program BEAST v.
434 1.8.3. [14] with the parameters used by Farrington et al. [20]: coalescent constant size prior to model
435 the branching process, strict clock, substitution rate equal to 1, specific models of DNA evolution de-
436 fined by the authors (supplementary table S2), and 10 million generations with sampling every 1000
437 generations. The monophyly of the ingroup and the topology among the outgroups were constrained
438 according to the topology depicted in [20].

439

440 *Estimation of parameters under the multi-species coalescent (MSC) model*

441 We optimized the MSC model for $N = 6$ species by varying two parameters, the speciation rate
442 λ and the extinction rate μ , and fixing the coalescence rate to 1. Birth-death trees of 6 tips (function
443 'sim.bdtree', package 'geiger' in R) were simulated in a grid of $(\lambda, \mu = m\lambda)$ with $\lambda \in [0.02, 0.34]$, every
444 0.02, and $m \in [0.1, 0.65]$, every 0.05. Because we simulated small trees (6 tips), the degree of variation
445 between trees simulated with the same parameters was high. Therefore for each value of (λ, μ) we
446 randomly selected 15 species trees for which the crown age did not differ by more than 2.5% from the
447 expected crown age. Next, we simulated 10 gene genealogies for each species tree (coalescence
448 rate fixed to 1).

449 If the diversification rate (speciation rate minus extinction rate) is low, all the homologous genes
450 will coalesce before the next node in the species tree, so that all the gene trees will have the same
451 topology. On the contrary, if the diversification rate is too fast, some homologous genes will not have
452 time to coalesce before the next node of the species tree, resulting in incongruent gene trees due to
453 the randomness of coalescences (ILS).

454

455 *Estimation of parameters under the gene-based diversification (GBD-backward) model*

456 Equivalently, we optimized the GBD-backward model for $N = 6$ by varying two parameters, here
457 a and b , and fixing $d = 1$ and $c = 200$ (recall c is given a sufficiently large value that coalescences
458 are instantaneous). Since increasing n has no effect on BHV distances (see results and figure 6), we
459 simulated genomes with $n = 10$ genes. The number of time units t was set to 5,000, which guaran-
460 tees the coalescence of all homologous genes. We performed 15 replicates under each parameter
461 combination in a grid of $(\frac{1}{a}, b)$ with $\frac{1}{a} \in [0.3, 3.5]$, every 0.2, and $b \in [0.01, 0.12]$, every 0.01.

462 For both models (MSC and GBD-backward) we employed the Kullback-Leibler (KL) divergence
463 (package 'FNN' in R) as a distance metric to find the best set of parameters by minimizing this dis-
464 tance between the distributions of BHV pairwise distances of empirical and simulated trees. The lower

465 the KL divergence is the better is the fit.

466

467 **RESULTS**

468 **Inference method accuracy**

469 Using simulated data-sets, we showed that our inference method was able to give reliable es-
470 timates of simulated parameters despite its simplicity (supplementary figure 1). Even if the exact
471 parameter combination was retrieved only twice over eight (test data sets d and e), the inferred pa-
472 rameters were very close to the simulated ones. We calculated the mean squared error (MSE),
473 defined as the average squared difference between the observed (inferred parameters) and predicted
474 values (simulated parameters). We found a MSE of $1.5e-04$ for the parameter $\frac{d}{a}$ and a MSE of 0.15
475 for the parameter b . This simple inference method is sufficient to estimate the parameters of the
476 model having supposedly shaped the gene trees of the data set. More subtle methods will be devel-
477 oped in the future to account for more complex features, such as differential gene flow depending on
478 putative gene categories, and to infer the very history of the embedding of gene lineages into species.

479

480 **Comparison of the GBD-models**

481 Even if the two models, GBD-forward and GBD-backward, are only approximately equal, they
482 showed a qualitatively similar pattern in gene genealogies and in dissimilarity among gene trees with
483 increasing gene flow (figure 3 and 4). Because they are co-adapted, genes sampled in the same
484 species at present time should have spent time together in the same species more often than by
485 chance in the past. This property was indeed observed in both models, with genes sampled at
486 present time frequently found together in the same species in the past (figure 3).

487 Using our inference method, we found a strong correlation between $\frac{d}{a}$ (GBD-backward) and $\frac{\delta}{\alpha}$
488 (GBD-forward) for $\beta = 0.05$ ($r = 0.99$ and p.value = 0.005) and $\beta = 0.02$ ($r = 0.97$ and p.value
489 = 0.03). For $\beta = 0.01$ we found a high correlation but not significant correlation (presumably due to
490 small sample size) ($r = 0.89$ and p.value = 0.1). Our inference method was unable to provide a good
491 estimate for β , this estimate oscillated between 0.01 and 0.02 regardless of β . However we found a
492 more pronounced slope for higher β indicating higher gene flow. If the colonization is fast, the number
493 of mutations differentially acquired within each population will be small, therefore introgression events
494 (gene flow) will be very likely among populations. With this inference method, the inclination of the
495 slope expresses the difference in β (colonization rate).

496

497 **A single sampled genome (GBD-backward)**

498 With $N = 1$ sampled genome containing n genes, we let $A(t) = (A_1(t), \dots, A_n(t))$ denote the
499 sorting of genes into ancestral species t units of time before the present. More precisely, $A_k(t)$
500 denotes the number of ancestral species containing k gene lineages, so that $n = \sum_{k=1}^n k A_k(t)$ and

501 $S(t) = \sum_{k=1}^n A_k(t)$ is the total number of species at t ancestral to the sampled genome. For each $\varepsilon \in$
 502 $(0, 1]$, we will also be interested in the number $S_\varepsilon(t) = \sum_{k=\lceil \varepsilon n \rceil}^n A_k(t)$ of ancestral species containing
 503 at least a fraction ε of the genome (with $\lceil x \rceil$ denoting the smallest integer larger than x). All stationary
 504 quantities will be denoted by the same symbols, replacing t with ∞ .

505 We will call a *block* at (backward) time t a (maximal) set of gene lineages that lie in the same
 506 species at time t . The transition rates can be specified as follows in terms of the configuration of
 507 gene lineages into blocks (*i.e.*, ancestral species). For each pair of blocks containing (j, k) lineages,
 508 non-homologous attraction occurs at rate ajk and results in the configuration $(j - 1, k + 1)$. For each
 509 block containing j lineages, gene flow occurs at rate dj and results in the block losing one lineage;
 510 simultaneously a new block containing 1 single lineage is created. These are exactly the same rates
 511 as in the well-known Moran model with mutation under the infinite-allele model [59], replacing ‘block’
 512 with ‘allele’, ‘connection’ by ‘resampling’ (simultaneous birth from one of the j carriers of a given allele
 513 and death of one of the k carriers of another given allele) and ‘gene flow’ with ‘mutation’ (mutation
 514 appearing in one of the j carriers of a given allele into a new allele never existing before). For this
 515 Moran model,

- 516 • the total population size is n ;
- 517 • at rate a for each oriented pair of individuals independently, the first individual of the pair gives
 518 birth to a copy of herself and the second individual of the pair is simultaneously killed;
- 519 • mutation occurs at rate d independently in each individual lineage.

As a consequence, $A(t)$ has the same distribution as the allele frequency spectrum in the Moran
 model with total population size n , resampling rate a and mutation rate d , starting at time $t = 0$
 from a population of clonal individuals (one single block). In particular, the distribution of $A(\infty)$ is
 the stationary distribution of the allele frequency spectrum, which is known to be given by Ewens’
 sampling formula with scaled mutation rate d/a [15, 18, 19]. Expectations of this distribution are:

$$\mathbb{E}(A_k(\infty)) = \frac{d}{d + a(k - 1)},$$

520 so that

$$\mathbb{E}(S(\infty)) = \sum_{k=1}^n \frac{d}{d + a(k - 1)} \tag{1}$$

521 and

$$\mathbb{E}(S_\varepsilon(\infty)) = \sum_{k=\lceil \varepsilon n \rceil}^n \frac{d}{d + a(k - 1)}. \tag{2}$$

In particular, as $n \rightarrow \infty$,

$$\mathbb{E}(S(\infty)) \sim \frac{d}{a} \ln(n) \quad \text{and} \quad \mathbb{E}(S_\varepsilon(\infty)) \sim \frac{d}{a} \ln(1/\varepsilon).$$

522

523 At stationarity, and particularly for large values of $\frac{d}{a}$, the mean number of ancestral species $S(\infty)$
524 obtained from simulations was equal to the mathematical prediction (figure 5A). In particular, the
525 mean number of ancestral species at stationarity increases with $\frac{d}{a}$.

526 An additional key feature of this model is *sampling consistency*. In words, the history of a sample
527 of k genes taken from a genome of n genes does not depend on n . This property can again be
528 deduced from the representation of our model in terms of the better known Moran model. Indeed,
529 the dynamics of a sample of k individuals in the Moran model does not depend on the population
530 size, as can be seen from the so-called lockdown construction [16]. The simulations performed with k
531 genes randomly sampled from each genome of n genes, are in agreement with this claim of sampling
532 consistency: the number of ancestral species at stationarity $\mathbb{E}(S_\varepsilon(\infty))$ is independent of the number
533 of genes n (figure 5B).

534

535 **A sample of several genomes (GBD-backward)**

536 Using simulations, we evaluated the GBD-backward model for several sampled genomes ($N >$
537 1) under several combinations of parameters. As expected gene tree diversity, measured by BHV
538 distances, increased with $\frac{d}{a}$, *i.e.* the relative amount of gene flow, and with the number of species N .
539 Conversely our results showed that the number of genes n had no effect on distances (figure 6A).
540 This last result, the lack of influence of n on gene tree diversity, is of particular interest, because one
541 usually has only access to a fraction of a genome. It shows that regardless of the number of genes
542 sampled, the resulting gene tree diversity will remain the same as long as gene trees have been
543 shaped by processes with similar parameter values.

544 Our results also showed that as the *homologous attraction* rate b decreases, and for the same $\frac{d}{a}$,
545 gene trees were more similar (lower BHV distances) (figure 6B). When a long period of time elapses
546 between two homologous attraction events (low b), all the genes belonging to the two genomes that
547 have started to coalesce, have enough time to converge toward the same species, and thus coalesce
548 before the next homologous attraction event, in spite of gene flow.

549

550 **GBD versus MSC: ignoring gene flow may lead to mistaken phylogenetic inferences**

551 When evaluating the ability of MSC model to deal with gene flow, we found a strong support (pos-
552 terior probabilities > 0.90) for all the nodes of the Bayesian species tree even if the individual gene
553 trees of the GBD-backward model did not corroborate this topology (figure 7). For example, 7 out
554 of 10 gene trees modeled under the GBD-backward model support the connection between species
555 E and species C and D, and only 3 the direct relationship between species E and F. **On the con-**
556 **trary**, the Bayesian tree strongly supports the clade (E,F) with a posterior probability equal to 1, and
557 considers all the connections between E and (C,D) to be due to ancestral polymorphism (*i.e.*, ILS).
558 Moreover because gene trees are constrained in the species tree (MSC model), the coalescences

559 between genes of E and (C,D) must take place after the species tree coalescence, therefore these
560 coalescences are timed around 7 My instead of 2 My according to the GBD-backward tree. Failing to
561 recognize that gene flow may have shaped gene genealogies, hence DNA sequences, can result in
562 important topological and dating errors.

563

564 **The GBD-backward model correctly captures the signal left by gene flow in empirical data-sets**

565 To find the best set of parameters, we minimized the Kullback-Leibler (KL) divergence between
566 the distributions of BHV pairwise distances of empirical and simulated trees (figure 8). Under the
567 multi-species coalescent (MSC) model, the most likely set of parameters was $\mu = 0.4 \times \lambda$ and $\lambda = 0.2$
568 (KL divergence = 0.23) for the bears, and $\mu = 0.45 \times \lambda$ and $\lambda = 0.22$ for the finches (KL divergence
569 = 0.12). We noted longer tailed distributions for the distances between trees modeled under the MSC
570 model than for the empirical data-sets (figure 9). This skewed distribution obtained with the MSC
571 model explains why we did not detect a sharp peak in the optimization landscape for the MSC model
572 (figure 8).

573 Under the gene-based diversification (GBD-backward) model, the most likely set of parameters
574 was $b = 0.03$ and $\frac{d}{a} = 2.1$ (KL divergence = 0.14) for the bears, and $b = 0.11$ and $\frac{d}{a} = 1.5$ for
575 the finches (KL divergence = 0.01) (figure 8). Contrary to the MSC model, the distributions of the
576 distances between trees modeled under the GBD-backward model or empirical trees did not show,
577 or to a lesser degree, a long tail (figure 9), explaining why we could detect a sharp peak in the
578 optimization landscape for the MSC model (figure 8).

579 Comparing the parameters λ and μ to b and $\frac{d}{a}$ is not straightforward as the two models, MSC
580 and GBD-backward, are built under different assumptions. However in both cases, the parameters
581 influence the diversity among trees (shape of the distribution of BHV pairwise distances). A greater
582 diversity among trees is expected with increasing λ and decreasing μ , and with increasing $\frac{d}{a}$ and b ,
583 allowing us to explore the parameter landscape to find the setting that minimizes the distance between
584 simulations and empirical data-sets for each model.

585 Given our results and the mathematical predictions, the time-averaged number $S_\varepsilon(\infty)$ of ancestral
586 species to the sampled genome containing at least 10% of the genome ($\varepsilon = 0.1$) when $n \rightarrow \infty$ is 4.8
587 for the bear data-set and 3.4 for the finch data-set.

588

589 **DISCUSSION**

590 Within species, gene flow allows the maintenance of species cohesion in the face of genetic
591 differentiation [60, 81], preventing genetic isolation of populations and the subsequent emergence of
592 reproductive barriers leading to speciation [10]. Among species, the existence of gene flow challenges
593 the notion of a species genealogy as well as the current concepts of species. Indeed, if gene flow is
594 as pervasive as recent empirical studies suggest [8, 11, 23, 37], the genealogical history of species

595 should be represented as a phylogenetic network encompassing the mosaic of gene genealogies.
596 Similarly, it seems very conservative to delineate species based on the widely used biological species
597 concept (reproductive isolation) [54], or phylogenetic species concept (reciprocal monophyly) [65].
598 Because of the ubiquity of gene flow, which can persist for several millions of years after the lineages
599 have started to diverge (*i.e.*, onset of speciation) [4, 49], species should be rather defined by their
600 capacity to coexist without fusion in spite of gene flow [50, 74].

601 The simplified view of diversification, consisting in representing lineages splitting instantaneously
602 into divergent lineages with no interaction (gene exchange) after the split, has been preventing evo-
603 lutionary biologists from fully apprehending diversification at the genomic level and from correctly
604 interpreting discrepancies between gene histories. Indeed, conflicting gene trees make the interpre-
605 tation of their evolutionary history difficult. However, we argue that phylogenetic incongruence among
606 gene trees should not be considered as a nuisance, but rather as a meaningful biological signal re-
607 vealing some features of the dynamics of genetic differentiation and of gene flow through time and
608 across clades. Current phylogenetic methods rely on the assumption that gene trees are constrained
609 within the species tree, and that gene flow occurs infrequently between species. For many data-sets
610 such as sequence alignments of genomes sampled from young clades, such methods could lead to
611 an evolutionary misinterpretation of gene trees, and in the worst case to species trees with high node
612 support while the gene trees had very different evolutionary histories (see figure 7). These obser-
613 vations urge for a change of paradigm, where gene flow is fully part of the diversification model. To
614 consider the ubiquity of gene flow across the Tree of Life [and its broad effect on genomes](#) described
615 by many recent studies, we have developed a new framework focusing on gene genealogies and
616 relaxing the constraints inherent to the MSC paradigm. [This framework is implemented in a math-](#)
617 [ematical model that we named the gene-based diversification \(GBD-forward\) model.](#) We have also
618 [developed a complementary version of this model, the GBD-backward model, speeding-up the simu-](#)
619 [lations thanks to a coalescent approach.](#)

620

621 **The GBD-backward model**

622 Under the GBD-backward model, gene genealogies are governed by four parameters correspond-
623 ing to four biological processes, *coalescence* (colonization), *non-homologous attraction* (mutation),
624 *homologous attraction* (reproductive isolation), and *gene flow* (introgression) (figure 2).

625 *Homologous attraction* corresponds to finding the most recent common ancestor of the two species
626 at the genomic level. The time spent between homologous attraction events depends crucially on the
627 (phylogenetic distance of the) species sampled at the present. *Gene flow* corresponds to the in-
628 trogression of genetic material from one species into another species. *Non-homologous attraction*
629 models the genetic differentiation (mutational process). The slower genes accumulate mutations and
630 differentiate, the more time can be spent by gene lineages in different species. Hence when genomes

631 differentiate slowly, the rate of non-homologous attraction is low.

632 Each of these parameters influences differently the resulting tree diversity, *i.e.* the distribution of
633 the BHV distances among trees, that we used here as a summary statistic. Instead of focusing on
634 the main phylogenetic signal alone as done by the current phylogenetic methods, the GBD-backward
635 model makes use of the whole signal encompassed by all gene trees.

636 Higher amount of *gene flow* and reduced time to untangle gene genealogies before the connection
637 of two other genomes (*homologous attraction*) increase the diversity among trees. Conversely, when
638 homologous genes coalesce faster (*coalescence*) and genes converge faster toward the species
639 harboring the other genes of their genome (*non-homologous attraction*) a lower diversity among trees
640 is expected.

641 After evaluating this model under various sets of parameters, we applied it to analyze two empiri-
642 cal multi-locus data-sets for which gene tree conflicts obscure the evolutionary history.

643

644 **Gene flow among bears and among finches**

645 Our results showed support for the hypothesis that gene flow has shaped the gene trees of bears
646 and finches (figure 9). For the bear data-set, we found that each species had on average in the past
647 about 4.8 ancestral species carrying at least 10% of its present genome (equation (2)). This result
648 is in line with previous studies reporting gene flow between pairs of bear species [7, 32, 42, 46, 56].
649 Moreover, a recent phylogenomic study (869 Mb divided into 18,621 genome fragments) confirmed
650 the existence of gene flow between sister species as well as between more phylogenetically distant
651 species [41]. The authors used the *D*-statistic (gene flow between sister species) and *D_{FOIL}*-
652 statistic (gene flow among ancestral lineages [66]) to detect gene flow among the 6 bear species. Using
653 their results, for each pair of species *ij* among the *N* species, we determined if the species *j* has
654 contributed ($g_{ij} = 1$) or not ($g_{ij} = 0$) to the genome of the species *i* (with $g_{ii} = 1$), and calculated the
655 average number of ancestral species *S* as follow:

$$\bar{S} = \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^N g_{ij}. \quad (3)$$

656 We found on average 5.3 ancestral species for each of the Ursinae bears [41], close to the estimate
657 obtained with the GBD-backward model (4.8).

658 We detected lower gene flow among finches than among bears. Each finch species had on aver-
659 age in the past 3.4 ancestral species (for the subsample of gene trees analyzed here), which is also
660 consistent with the extensive evidence that many species hybridize on several islands [22, 28, 30, 31,
661 75]. Because of gene flow very little genetic structure was detected by a Bayesian population struc-
662 ture analysis, only 3 genetic populations among the 6 *Geospiza* species [20]. Each of the 2 species,
663 *G. magnirostris* and *G. scandens*, were mostly characterized by a single genetic population, there-
664 fore had about 1 ancestral species each. Conversely 4 *Geospiza* species shared the same genetic

665 population, suggesting 4 ancestral species for each of these 4 species. Taking together these results
666 roughly indicate that each of the 6 *Geospiza* species had in average 3 ancestral species, in line with
667 the GBD-backward estimate (3.4).

668 In many cases, such as among bears and finches, gene flow is frequent and complicates the
669 relationships between species, challenging the notion of a unique species tree. A strictly bifurcating
670 lineage-based model will not adequately reflect those complex evolutionary patterns. On the contrary,
671 models developed under the *genomic view of diversification* framework, *i.e.* relaxing species bound-
672 aries and accounting for gene flow, will better reproduce the complex history of gene genealogies
673 under pervasive gene flow. Note that we considered a simple scenario with no ILS and statistically
674 exchangeable genes resulting in a model with only three parameters, but given the simplicity and
675 the flexibility of our model, many extensions may be considered to address scenarios that could not
676 have been considered previously, opening up new perspectives in the study of speciation and macro-
677 evolution.

678

679 **Gene flow: an evolutionary force driving diversification**

680 Species diversification requires genetic variation among organisms, introduced by mutations and
681 structural variation, upon which natural selection and drift can act by influencing the sorting of offspring
682 and the survival of organisms [74]. Recently, gene flow has also been mentioned as another potential
683 source of genetic variation [53], and more particularly in the case of adaptive radiations [9, 44, 55, 78].
684 Hybrid zones act as filters, preventing the introgression of deleterious genes while allowing advanta-
685 geous or neutral genes to cross the species boundaries [53]. Newly acquired genes will then be a
686 source of variation [53], by providing evolutionary adaptive shortcuts or greater adaptability once in
687 the genetic pool of the introgressed species [53]. The introgressed species then has a wider range of
688 potentially adaptive allelic variants, allowing it to diversify rapidly if the opportunity arises. Accordingly
689 important gene flow should be detected prior to an adaptive radiation. This hypothesis is supported
690 by empirical evidence, but has only been tested under limited conditions [9, 44, 55, 78]. The model
691 proposed here constitutes an opportunity to investigate more systematically how gene flow is dis-
692 tributed throughout the phylogenies and [if gene flow can facilitate adaptive radiations](#).

693

694 **Evolutionary dynamics along the genome**

695 Along the genome, gene flow is not expected to be uniformly distributed either. Incongruent gene
696 trees can make genes that have evolved more slowly stand out. [Indeed, gene flow among populations](#)
697 [undergoing divergent selection will depend on the number of new alleles acquired differentially \(non](#)
698 [co-adapted\) within each population. Through the action of introgression and recombination, gene flow](#)
699 [will persist longer among genomic regions undergoing a slow genetic differentiation](#). Conversely, con-
700 congruent gene trees should reveal genomic regions not subject to gene flow, like genomic regions under

701 strong selective differentiation, *i.e.*, regions that harbor rapidly fixed divergent (non co-adapted) alle-
702 les [33, 36]. This framework could thus be used to evaluate how gene flow varies along the genome
703 and to explore the genomic architecture of species barriers. Indeed some regions, as sexual chromo-
704 somes or low recombination regions, are expected to be more differentiated and hence to undergo
705 less gene flow (e.g. *Heliconius* species [52]). In order to distinguish between genes and to reduce
706 potential errors in parameter estimation, data may be grouped by gene class (statistical binning) using
707 a method aiming to evaluate whether two genes are likely to have the same tree (linkage) or the same
708 tree in distribution (statistical exchangeability) [57].

709

710 Perspectives

711 Phylogenetic models and methods inferring macro-evolutionary history, such as speciation and
712 extinction rates, trait evolution or ancestral characters, have become increasingly complex [61, 70, 85].
713 Yet, the raw material used by these methods is often reduced to the species tree, which can be viewed
714 as a summary statistic of the information contained in the genome. We argue here that a valuable
715 amount of additional signal, not accessible in species trees, is contained in gene trees, and is directly
716 informative about the diversification process. Indeed, because genetic differentiation and gene flow
717 impact each gene differently, genes may have experienced very different evolutionary trajectories.

718 In order to make use of the entire information conveyed by gene trees, we have proposed here
719 a new approach to study diversification, the genomic view of diversification, under which gene trees
720 shape the species tree rather than the opposite. This approach aims at better depicting the intri-
721 cate evolutionary history of species and genomes. We hope that this view of diversification will pave
722 the way for future developments in the perspective of inferring diversification processes directly from
723 genomes rather than from their summary into one single species tree. One of the challenges in
724 this direction will be to propose finer inference methods than the simple, but reasonably satisfactory,
725 method used here, based on a single multidimensional summary statistic, the distribution of pairwise
726 BHV distances between gene trees.

727

728 SUPPLEMENTARY MATERIAL

729 Supplementary Material and code for the models are deposited on bioRxiv.

730

731 ACKNOWLEDGMENTS

732 The authors thank the *Center for Interdisciplinary Research in Biology* (Collège de France, CNRS)
733 for funding. JM is funded by LabEx MemoLife, project *Genomics of Diversification*. The authors also
734 thank the INRA MIGALE bioinformatics platform (<http://migale.jouy.inra.fr>) for providing computational
735 resources. The authors warmly thank the Recommender of *PCI Evolutionary Biology* (Peter Ralph)
736 and two anonymous reviewers for their very constructive and relevant comments on the first version

737 of this paper.

738

739 References

- 740 [1] ABBOTT, R., ALBACH, D., ANSELL, S., ARNTZEN, J. W., BAIRD, S. J. E., BIERNE, N., BOUGH-
741 MAN, J., BRELSFORD, A., BUERKLE, C. A., BUGGS, R., BUTLIN, R. K., DIECKMANN, U., ER-
742 OUKHMANOFF, F., GRILL, A., CAHAN, S. H., HERMANSEN, J. S., HEWITT, G., HUDSON, A. G.,
743 JIGGINS, C., JONES, J., KELLER, B., MARCZEWSKI, T., MALLET, J., MARTINEZ-RODRIGUEZ, P.,
744 MÖST, M., MULLEN, S., NICHOLS, R., NOLTE, A. W., PARISOD, C., PFENNIG, K., RICE, A. M.,
745 RITCHIE, M. G., SEIFERT, B., SMADJA, C. M., STELKENS, R., SZYMURA, J. M., VÄINÖLÄ, R.,
746 WOLF, J. B. W., AND ZINNER, D. Hybridization and speciation. *Journal of Evolutionary Biology*
747 *26*, 2 (2013), 229–246.
- 748 [2] BERLOCHER, S. H. Radiation and divergence in the *Rhagoletis pomonella* species group: infer-
749 ences from allozymes. *Evolution* *54*, 2 (2000), 543–557.
- 750 [3] BILLERA, L. J., HOLMES, S. P., AND VOGTMANN, K. Geometry of the space of phylogenetic
751 trees. *Advances in Applied Mathematics* *27*, 4 (2001), 733–767.
- 752 [4] BOLNICK, D. I., NEAR, T. J., AND NOOR, M. Tempo of hybrid inviability in centrarchid fishes
753 (teleostei: centrarchidae). *Evolution* *59*, 8 (2005), 1754–1767.
- 754 [5] BOUCKAERT, R., HELED, J., KÜHNERT, D., VAUGHAN, T., WU, C.-H., XIE, D., SUCHARD, M. A.,
755 RAMBAUT, A., AND DRUMMOND, A. J. BEAST 2: A Software Platform for Bayesian Evolutionary
756 Analysis. *PLOS Computational Biology* *10*, 4 (Apr. 2014), e1003537.
- 757 [6] BUONERBA, L., ZACCARA, S., DELMASTRO, G. B., LORENZONI, M., SALZBURGER, W., AND
758 GANTE, H. F. Intrinsic and extrinsic factors act at different spatial and temporal scales to shape
759 population structure, distribution and speciation in Italian *Barbus* (Osteichthyes: Cyprinidae).
760 *Molecular Phylogenetics and Evolution* *89* (2015), 115–129.
- 761 [7] CAHILL, J. A., GREEN, R. E., FULTON, T. L., STILLER, M., JAY, F., OVSYANIKOV, N.,
762 SALAMZADE, R., JOHN, J. S., STIRLING, I., AND SLATKIN, M. Genomic evidence for island
763 population conversion resolves conflicting theories of polar bear evolution. *PLoS genetics* *9*, 3
764 (2013), e1003345.
- 765 [8] CLARK, A. G., AND MESSER, P. W. Conundrum of jumbled mosquito genomes. *Science* *347*,
766 6217 (2015), 27–28.

- 767 [9] CONSORTIUM, H. G. Butterfly genome reveals promiscuous exchange of mimicry adaptations
768 among species. *Nature* 487, 7405 (2012), 94–98.
- 769 [10] COYNE, J. A., AND ORR, H. A. *Speciation*. Sinauer Associates, Sunderland, MA, 2004.
- 770 [11] CUI, R., SCHUMER, M., KRUESI, K., WALTER, R., ANDOLFATTO, P., AND ROSENTHAL, G. G.
771 Phylogenomics reveals extensive reticulate evolution in Xiphophorus fishes. *Evolution* 67, 8
772 (2013), 2166–2179.
- 773 [12] DE BUSSCHERE, C., HENDRICKX, F., VAN BELLEGHEM, S. M., BACKELJAU, T., LENS, L., AND
774 BAERT, L. Parallel habitat specialization within the wolf spider genus Hogna from the Galápagos.
775 *Molecular ecology* 19, 18 (2010), 4029–4045.
- 776 [13] DOBZHANSKY, T. H. Studies on hybrid sterility. II. Localization of sterility factors in *Drosophila*
777 *pseudoobscura* hybrids. *Genetics* 21, 2 (1936), 113.
- 778 [14] DRUMMOND, A. J., SUCHARD, M. A., XIE, D., AND RAMBAUT, A. Bayesian phylogenetics with
779 BEAUti and the BEAST 1.7. *Molecular biology and evolution* 29, 8 (2012), 1969–1973.
- 780 [15] DURRETT, R. *Probability Models for DNA Sequence Evolution*. Springer, Dec. 2008. Google-
781 Books-ID: o4_bMHY7jFoC.
- 782 [16] ETHERIDGE, A. *Some Mathematical Models from Population Genetics: École D'Été de Prob-*
783 *abilités de Saint-Flour XXXIX-2009*. Springer Science & Business Media, Jan. 2011. Google-
784 Books-ID: miI9tdPCFdUC.
- 785 [17] ETIENNE, R. S., MORLON, H., AND LAMBERT, A. Estimating the duration of speciation from
786 phylogenies. *Evolution* 68, 8 (2014), 2430–2440.
- 787 [18] EWENS, W. J. The sampling theory of selectively neutral alleles. *Theoretical Population Biology*
788 3, 1 (Mar. 1972), 87–112.
- 789 [19] EWENS, W. J., AND TAVARÉ, S. Ewens Sampling Formula. In *Encyclopedia of Statistical Sci-*
790 *ences*. American Cancer Society, 2006.
- 791 [20] FARRINGTON, H. L., LAWSON, L. P., CLARK, C. M., AND PETREN, K. The evolutionary history of
792 Darwin's finches: speciation, gene flow, and introgression in a fragmented landscape. *Evolution*
793 68, 10 (2014), 2932–2944.
- 794 [21] FONTENOT, B. E., MAKOWSKY, R., AND CHIPPINDALE, P. T. Nuclear–mitochondrial discordance
795 and gene flow in a recent radiation of toads. *Molecular Phylogenetics and Evolution* 59, 1 (Apr.
796 2011), 66–80.

- 797 [22] FREELAND, J. R., AND BOAG, P. T. THE MITOCHONDRIAL AND NUCLEAR GENETIC HOMO-
798 GENEITY OF THE PHENOTYPICALLY DIVERSE DARWIN'S GROUND FINCHES. *Evolution;*
799 *International Journal of Organic Evolution* 53, 5 (Oct. 1999), 1553–1563.
- 800 [23] GALLUS, S., JANKE, A., KUMAR, V., AND NILSSON, M. A. Disentangling the relationship of the
801 Australian marsupial orders using retrotransposon and evolutionary network analyses. *Genome*
802 *biology and evolution* 7, 4 (2015), 985–992.
- 803 [24] GANTE, H. F., COLLARES-PEREIRA, M. J., AND COELHO, M. M. Introgressive hybridisation
804 between two Iberian Chondrostoma species (Teleostei, Cyprinidae) revisited: new evidence from
805 morphology, mitochondrial DNA, allozymes and NOR-phenotypes. *Folia Zoologica* 53, 4 (2004),
806 423.
- 807 [25] GANTE, H. F., DOADRIO, I., ALVES, M. J., AND DOWLING, T. E. Semi-permeable species
808 boundaries in Iberian barbels (Barbus and Luciobarbus, Cyprinidae). *BMC evolutionary biology*
809 15, 1 (2015), 111.
- 810 [26] GANTE, H. F., SANTOS, C. D., AND ALVES, M. J. Phylogenetic relationships of the newly
811 described species Chondrostoma olisiponensis (Teleostei: Cyprinidae). *Journal of Fish Biology*
812 76, 4 (Mar. 2010), 965–974.
- 813 [27] GRANT, B. R., AND GRANT, P. R. Hybridization and speciation in Darwin's finches: the role of
814 sexual imprinting on a culturally transmitted trait. *Endless forms: species and speciation* (1998),
815 404–422.
- 816 [28] GRANT, P. R., AND GRANT, B. R. PHENOTYPIC AND GENETIC EFFECTS OF HYBRIDIZA-
817 TION IN DARWIN'S FINCHES. *Evolution; International Journal of Organic Evolution* 48, 2 (Apr.
818 1994), 297–316.
- 819 [29] GRANT, P. R., AND GRANT, B. R. Speciation and hybridization in island birds. *Phil. Trans. R.*
820 *Soc. Lond. B* 351, 1341 (1996), 765–772.
- 821 [30] GRANT, P. R., AND GRANT, B. R. Hybridization, Sexual Imprinting, and Mate Choice. *The*
822 *American Naturalist* 149, 1 (1997), 1–28.
- 823 [31] GRANT, P. R., GRANT, B. R., AND PETREN, K. Hybridization in the recent past. *The American*
824 *Naturalist* 166, 1 (July 2005), 56–67.
- 825 [32] HAILER, F., KUTSCHERA, V. E., HALLSTRÖM, B. M., KLASSERT, D., FAIN, S. R., LEONARD,
826 J. A., ARNASON, U., AND JANKE, A. Nuclear Genomic Sequences Reveal that Polar Bears Are
827 an Old and Distinct Bear Lineage. *Science* 336, 6079 (Apr. 2012), 344–347.

- 828 [33] HARRISON, R. G., AND LARSON, E. L. Hybridization, introgression, and the nature of species
829 boundaries. *Journal of Heredity* 105, S1 (2014), 795–809.
- 830 [34] HEDGES, S. B., AND KUMAR, S. *The Timetree of Life*. OUP Oxford, Apr. 2009. Google-Books-
831 ID: 9rt1c1hl49MC.
- 832 [35] HELED, J., AND DRUMMOND, A. J. Bayesian inference of species trees from multilocus data.
833 *Molecular Biology and Evolution* 27, 3 (2010), 570–580.
- 834 [36] JANOUŠEK, V., MUNCLINGER, P., WANG, L., TEETER, K. C., AND TUCKER, P. K. Functional
835 organization of the genome may shape the species boundary in the house mouse. *Molecular*
836 *biology and evolution* 32, 5 (2015), 1208–1220.
- 837 [37] JÓNSSON, H., SCHUBERT, M., SEGUIN-ORLANDO, A., GINOLHAC, A., PETERSEN, L., FUMA-
838 GALLI, M., ALBRECHTSEN, A., PETERSEN, B., KORNELIUSSEN, T. S., VILSTRUP, J. T., LEAR,
839 T., MYKA, J. L., LUNDQUIST, J., MILLER, D. C., ALFARHAN, A. H., ALQURAISHI, S. A., AL-
840 RASHEID, K. A. S., STAGEGAARD, J., STRAUSS, G., BERTELSEN, M. F., SICHERITZ-PONTEN,
841 T., ANTCZAK, D. F., BAILEY, E., NIELSEN, R., WILLERSLEV, E., AND ORLANDO, L. Speciation
842 with gene flow in equids despite extensive chromosomal plasticity. *Proceedings of the National*
843 *Academy of Sciences* 111, 52 (2014), 18655–18660.
- 844 [38] KISHINO, H., THORNE, J. L., AND BRUNO, W. J. Performance of a Divergence Time Estimation
845 Method under a Probabilistic Model of Rate Evolution. *Molecular Biology and Evolution* 18, 3
846 (Mar. 2001), 352–361.
- 847 [39] KNOWLES, L. L., AND KUBATKO, L. S. *Estimating species trees: practical and theoretical as-*
848 *pects*. John Wiley and Sons, 2011.
- 849 [40] KUBATKO, L. S. Identifying hybridization events in the presence of coalescence via model selec-
850 tion. *Systematic Biology* 58, 5 (2009), 478–488.
- 851 [41] KUMAR, V., LAMMERS, F., BIDON, T., PFENNINGER, M., KOLTER, L., NILSSON, M. A., AND
852 JANKE, A. The evolutionary history of bears is characterized by gene flow across species.
853 *Scientific Reports* 7 (Apr. 2017), 46487.
- 854 [42] KUTSCHERA, V. E., BIDON, T., HAILER, F., RODI, J. L., FAIN, S. R., AND JANKE, A. Bears in
855 a forest of gene trees: phylogenetic inference is complicated by incomplete lineage sorting and
856 gene flow. *Molecular Biology and Evolution* 31, 8 (2014), 2004–2017.
- 857 [43] LAMBERT, A., MORLON, H., AND ETIENNE, R. S. The reconstructed tree in the lineage-based
858 model of protracted speciation. *Journal of mathematical biology* 70, 1-2 (2015), 367–397.

- 859 [44] LAMICHHANEY, S., BERGLUND, J., ALMÉN, M. S., MAQBOOL, K., GRABHERR, M., MARTINEZ-
860 BARRIO, A., PROMEROVÁ, M., RUBIN, C.-J., WANG, C., ZAMANI, N., GRANT, B. R., GRANT,
861 P. R., WEBSTER, M. T., AND ANDERSSON, L. Evolution of Darwin's finches and their beaks
862 revealed by genome sequencing. *Nature* 518, 7539 (2015), 371.
- 863 [45] LIU, L., YU, L., AND EDWARDS, S. V. A maximum pseudo-likelihood approach for estimating
864 species trees under the coalescent model. *BMC Evolutionary Biology* 10 (2010), 302.
- 865 [46] LIU, S., LORENZEN, E., FUMAGALLI, M., LI, B., HARRIS, K., XIONG, Z., ZHOU, L., KOR-
866 NELIUSSEN, T., SOMEL, M., BABBITT, C., WRAY, G., LI, J., HE, W., WANG, Z., FU, W., XIANG,
867 X., MORGAN, C., DOHERTY, A., O'CONNELL, M., MCINERNEY, J., BORN, E., DALÉN, L., DI-
868 ETZ, R., ORLANDO, L., SONNE, C., ZHANG, G., NIELSEN, R., WILLERSLEV, E., AND WANG,
869 J. Population Genomics Reveal Recent Speciation and Rapid Evolutionary Adaptation in Polar
870 Bears. *Cell* 157, 4 (May 2014), 785–794.
- 871 [47] LONG, C., AND KUBATKO, L. The effect of gene flow on coalescent-based species-tree inference.
872 *generations* 1 (2018), 2N.
- 873 [48] MADDISON, W. P. Gene trees in species trees. *Systematic Biology* 46, 3 (1997), 523–536.
- 874 [49] MALLETT, J. Hybridization as an invasion of the genome. *Trends in ecology & evolution* 20, 5
875 (2005), 229–237.
- 876 [50] MALLETT, J. Hybridization, ecological races and the nature of species: empirical evidence for the
877 ease of speciation. *Philosophical Transactions of the Royal Society B: Biological Sciences* 363,
878 1506 (2008), 2971–2986.
- 879 [51] MALLETT, J., BESANSKY, N., AND HAHN, M. W. How reticulated are species? *BioEssays* 38, 2
880 (2016), 140–149.
- 881 [52] MARTIN, S. H., DASMAHAPATRA, K. K., NADEAU, N. J., SALAZAR, C., WALTERS, J. R., SIMP-
882 SON, F., BLAXTER, M., MANICA, A., MALLETT, J., AND JIGGINS, C. D. Genome-wide evidence
883 for speciation with gene flow in *Heliconius* butterflies. *Genome Research* 23, 11 (Jan. 2013),
884 1817–1828.
- 885 [53] MARTINSEN, G. D., WHITHAM, T. G., TUREK, R. J., AND KEIM, P. Hybrid populations selectively
886 filter gene introgression between species. *Evolution* 55, 7 (2001), 1325–1335.
- 887 [54] MAYR, E. *Systematics and the origin of species, from the viewpoint of a zoologist*. Harvard
888 University Press, 1942.
- 889 [55] MEIER, J. I., MARQUES, D. A., MWAIKO, S., WAGNER, C. E., EXCOFFIER, L., AND SEEHAUSEN,
890 O. Ancient hybridization fuels rapid cichlid fish adaptive radiations. *Nature Communications* 8
891 (2017).

- 892 [56] MILLER, W., SCHUSTER, S. C., WELCH, A. J., RATAN, A., BEDOYA-REINA, O. C., ZHAO, F.,
893 KIM, H. L., BURHANS, R. C., DRAUTZ, D. I., AND WITTEKINDT, N. E. Polar and brown bear
894 genomes reveal ancient admixture and demographic footprints of past climate change. *Proceed-*
895 *ings of the National Academy of Sciences* 109, 36 (2012), E2382–E2390.
- 896 [57] MIRARAB, S., BAYZID, M. S., BOUSSAU, B., AND WARNOW, T. Statistical binning enables an
897 accurate coalescent-based estimation of the avian tree. *Science* 346, 6215 (2014), 1250463.
- 898 [58] MIRARAB, S., REAZ, R., BAYZID, M. S., ZIMMERMANN, T., SWENSON, M. S., AND WARNOW, T.
899 ASTRAL: genome-scale coalescent-based species tree estimation. *Bioinformatics* 30, 17 (Sept.
900 2014), i541–i548.
- 901 [59] MORAN, P. A. P. Random processes in genetics. In *Mathematical Proceedings of the Cambridge*
902 *Philosophical Society* (1958), vol. 54, Cambridge University Press, pp. 60–71.
- 903 [60] MORJAN, C. L., AND RIESEBERG, L. H. How species evolve collectively: implications of gene
904 flow and selection for the spread of advantageous alleles. *Molecular ecology* 13, 6 (2004), 1341–
905 1356.
- 906 [61] MORLON, H. Phylogenetic approaches for studying diversification. *Ecology Letters* 17, 4 (Apr.
907 2014), 508–525.
- 908 [62] MULLER, H. J. Recessive genes causing interspecific sterility and other disharmonies between
909 *Drosophila melanogaster* and *simulans*. *Genetics* 27 (1942), 157.
- 910 [63] NADEAU, N. J., MARTIN, S. H., KOZAK, K. M., SALAZAR, C., DASMAHAPATRA, K. K., DAVEY,
911 J. W., BAXTER, S. W., BLAXTER, M. L., MALLET, J., AND JIGGINS, C. D. Genome-wide patterns
912 of divergence and gene flow across a butterfly radiation. *Molecular Ecology* 22, 3 (Feb. 2013),
913 814–826.
- 914 [64] ORR, H. A. The population genetics of speciation: the evolution of hybrid incompatibilities.
915 *Genetics* 139, 4 (1995), 1805–1813.
- 916 [65] PAPADOPOULOU, A., BERGSTEN, J., FUJISAWA, T., MONAGHAN, M. T., BARRACLOUGH, T. G.,
917 AND VOGLER, A. P. Speciation and DNA barcodes: testing the effects of dispersal on the forma-
918 tion of discrete sequence clusters. *Philosophical Transactions of the Royal Society B: Biological*
919 *Sciences* 363, 1506 (2008), 2987–2996.
- 920 [66] PEASE, J. B., AND HAHN, M. W. Detection and Polarization of Introgression in a Five-Taxon
921 Phylogeny. *Systematic Biology* 64, 4 (July 2015), 651–662.
- 922 [67] PECCOUD, J., OLLIVIER, A., PLANTEGENEST, M., AND SIMON, J.-C. A continuum of genetic
923 divergence from sympatric host races to species in the pea aphid complex. *Proceedings of the*
924 *National Academy of Sciences* 106, 18 (May 2009), 7495–7500.

- 925 [68] PEREIRA, R. J., MONAHAN, W. B., AND WAKE, D. B. Predictors for reproductive isolation in a
926 ring species complex following genetic and ecological divergence. *BMC Evolutionary Biology* 11
927 (July 2011), 194.
- 928 [69] PINHO, C., AND HEY, J. Divergence with Gene Flow: Models and Data. *Annual Review of*
929 *Ecology, Evolution, and Systematics* 41, 1 (2010), 215–230.
- 930 [70] PYRON, R. A., AND BURBRINK, F. T. Phylogenetic estimates of speciation and extinction rates
931 for testing ecological and evolutionary hypotheses. *Trends in Ecology & Evolution* 28, 12 (Dec.
932 2013), 729–736.
- 933 [71] ROBINSON, D. F., AND FOULDS, L. R. Comparison of phylogenetic trees. *Mathematical bio-*
934 *sciences* 53, 1-2 (1981), 131–147.
- 935 [72] ROSINDELL, J., CORNELL, S. J., HUBBELL, S. P., AND ETIENNE, R. S. Protracted speciation
936 revitalizes the neutral theory of biodiversity. *Ecology Letters* 13, 6 (2010), 716–727.
- 937 [73] RUNDLE, H. D., NAGEL, L., BOUGHMAN, J. W., AND SCHLUTER, D. Natural selection and
938 parallel speciation in sympatric sticklebacks. *Science* 287, 5451 (2000), 306–308.
- 939 [74] SAMADI, S., AND BARBEROUSSE, A. The tree, the network, and the species. *Biological Journal*
940 *of the Linnean Society* 89, 3 (2006), 509–521.
- 941 [75] SATO, A., O'HUIGIN, C., FIGUEROA, F., GRANT, P. R., GRANT, B. R., TICHY, H., AND KLEIN, J.
942 Phylogeny of Darwin's finches as revealed by mtDNA sequences. *Proceedings of the National*
943 *Academy of Sciences* 96, 9 (Apr. 1999), 5101–5106.
- 944 [76] SCHLIEP, K. P. phangorn: phylogenetic analysis in R. *Bioinformatics* 27, 4 (2011), 592.
- 945 [77] SCHLUTER, D. Ecological causes of speciation. *Endless forms: species and speciation* (1998),
946 114–129.
- 947 [78] SEEHAUSEN, O. African cichlid fish: a model system in adaptive radiation research. *Proceedings*
948 *of the Royal Society of London B: Biological Sciences* 273, 1597 (2006), 1987–1998.
- 949 [79] SEEHAUSEN, O., BUTLIN, R. K., KELLER, I., WAGNER, C. E., BOUGHMAN, J. W., HOHENLOHE,
950 P. A., PEICHEL, C. L., AND SAETRE, G.-P. Genomics and the origin of species. *Nature Reviews*
951 *Genetics* 15, 3 (2014), 176–193.
- 952 [80] SIPOS, B., MASSINGHAM, T., JORDAN, G. E., AND GOLDMAN, N. PhyloSim-Monte Carlo simu-
953 lation of sequence evolution in the R statistical computing environment. *BMC bioinformatics* 12,
954 1 (2011), 104.

- 955 [81] SLATKIN, M. Gene flow and the geographic structure of natural populations. *Sci-*
956 *ence(Washington)* 236, 4803 (1987), 787–792.
- 957 [82] SOLÍS-LEMUS, C., YANG, M., AND ANÉ, C. Inconsistency of species tree methods under gene
958 flow. *Systematic Biology* 65, 5 (2016), 843–851.
- 959 [83] SOUCY, S. M., HUANG, J., AND GOGARTEN, J. P. Horizontal gene transfer: building the web of
960 life. *Nature Reviews Genetics* 16, 8 (Aug. 2015), 472–482.
- 961 [84] SOUSA-SANTOS, C., GANTE, H. F., ROBALO, J., CUNHA, P. P., MARTINS, A., ARRUDA, M.,
962 ALVES, M. J., AND ALMADA, V. Evolutionary history and population genetics of a cyprinid fish
963 (<Emphasis Type="Italic">Iberochondrostoma olisiponensis</Emphasis>) endangered by intro-
964 gression from a more abundant relative. *Conservation Genetics* 15, 3 (June 2014), 665–677.
- 965 [85] STADLER, T. Recovering speciation and extinction dynamics based on phylogenies. *Journal of*
966 *Evolutionary Biology* 26, 6, 1203–1219.
- 967 [86] STAMATAKIS, A. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large
968 phylogenies. *Bioinformatics* 30, 9 (2014), 1312–1313.
- 969 [87] TAMURA, K., BATTISTUZZI, F. U., BILLING-ROSS, P., MURILLO, O., FILIPSKI, A., AND KUMAR,
970 S. Estimating divergence times in large molecular phylogenies. *Proceedings of the National*
971 *Academy of Sciences* 109, 47 (Nov. 2012), 19333–19338.
- 972 [88] TURELLI, M., AND ORR, H. A. Dominance, epistasis and the genetics of postzygotic isolation.
973 *Genetics* 154, 4 (2000), 1663–1679.
- 974 [89] WAHLBERG, N., WEINGARTNER, E., WARREN, A. D., AND NYLIN, S. Timing major con-
975 flict between mitochondrial and nuclear genes in species relationships of Polygoni butterflies
976 (Nymphalidae: Nymphalini). *BMC Evolutionary Biology* 9 (May 2009), 92.
- 977 [90] WILLIS, S. C., MACRANDER, J., FARIAS, I. P., AND ORTÍ, G. Simultaneous delimitation of
978 species and quantification of interspecific hybridization in Amazonian peacock cichlids (genus
979 cichla) using multi-locus data. *BMC Evolutionary Biology* 12 (2012), 96.
- 980 [91] WU, C.-I. The genic view of the process of speciation. *Journal of Evolutionary Biology* 14, 6
981 (2001), 851–865.
- 982 [92] WU, Y. Coalescent-based species tree inference from gene tree topologies under incomplete
983 lineage sorting by maximum likelihood. *Evolution; international journal of organic evolution* 66, 3
984 (2012), 763–775.
- 985 [93] XU, B., AND YANG, Z. Challenges in Species Tree Estimation Under the Multispecies Coalescent
986 Model. *Genetics* 204, 4 (Dec. 2016), 1353–1368.

- 987 [94] YANG, Z. The BPP program for species tree estimation and species delimitation. *Current Zoology*
988 *61*, 5 (2015), 854–865.
- 989 [95] YU, Y., DONG, J., LIU, K. J., AND NAKHLEH, L. Maximum likelihood inference of reticulate
990 evolutionary histories. *Proceedings of the National Academy of Sciences* *111*, 46 (2014), 16448–
991 16453.

Figure 1: Gene trees and species tree conflicts. The species tree of A, B, and C is depicted in black. In pink (gene 1) and green (gene 2) are two gene trees congruent with the species tree, *i.e.* with A and B being sister species. In light blue (gene 3), the tree of a gene undergoing gene flow between species B and C. In dark blue (gene 4), the tree of a gene undergoing incomplete lineage sorting.

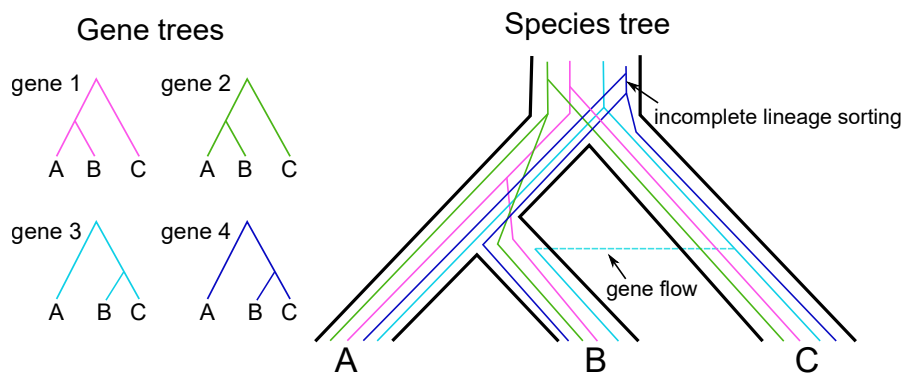
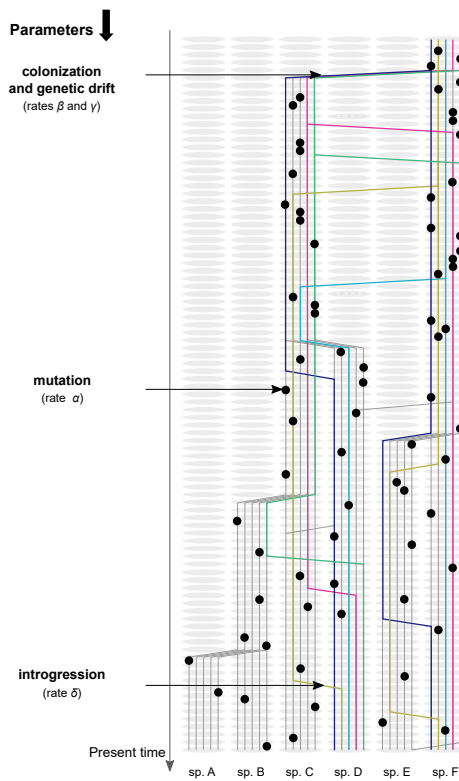


Figure 2: The gene-based diversification (GBD) models. Gene genealogies through species (or populations, depending on the point of view, retrospective vs prospective) are depicted for two present-day genomes ($N = 2$ at $t = 0$) and five homologous genes ($n = 5$). Each grey ellipse represents a species (A-F). The grey lines represent the gene genealogies of non-sampled species at $t = 0$. The model assumes that species are quasi-static in the timescale of a few generations, and each species lineage is located in a separate column. The genealogies of genes depend on four processes: introgression (*gene flow*), mutation (*non-homologous attraction*), colonization (*homologous attraction*), and genetic drift (*coalescence*).

GBD-forward model



GBD-backward model

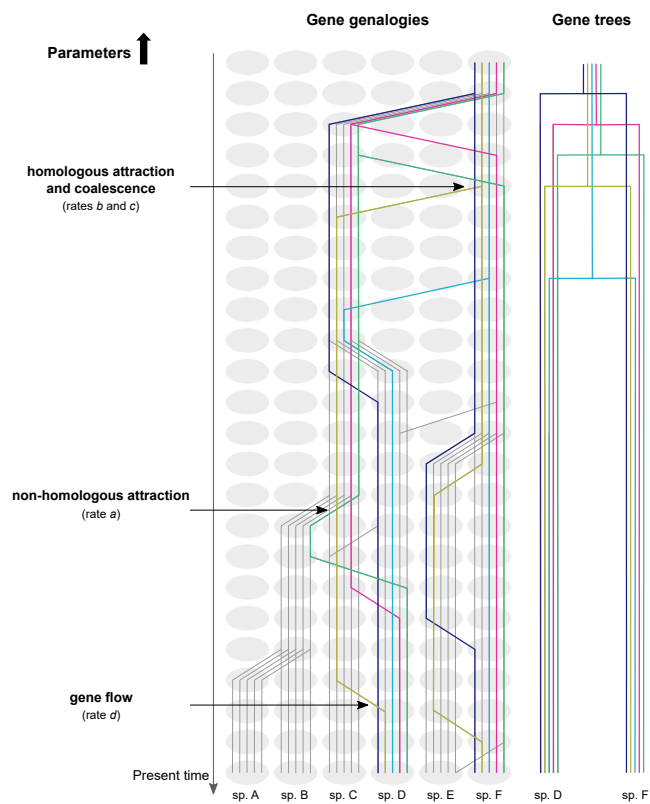


Figure 3: Genealogies of a single genome generated with the GBD-forward (A) and GBD-backward models (B). The labels/locations of species (or populations, depending on the point of view, retrospective vs prospective) are neutral. A) Parameter settings: $\alpha = 0.5$, $\beta = 1$, $\delta = 0.2$, $n = 5$ and $N = 30$. B) Parameter settings: $a = 1$, $b = 0.1$, $d = 2$, $n = 5$ and $N = 10$.

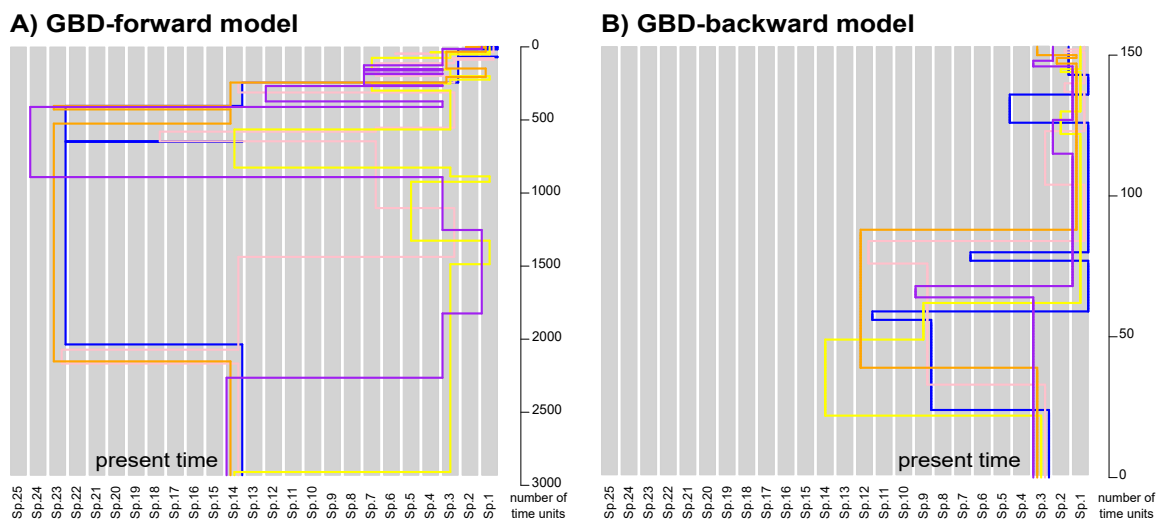


Figure 4: Comparison of the parameters of the GBD models. The Kullback-Leibler (KL) divergence was minimized between the distributions of BHV pairwise distances of GBD-forward and GBD-backward trees (with $N = 6$ and $n = 10$). The GBD-forward model was stopped when the number of populations was 30. Trees were built from the first 6 genomes (populations). Parameter settings: $\alpha = 0.5$, $\beta = 0.01, 0.02$ and 0.05 , $\delta \in [0.01, 0.04]$, every 0.01, $n = 10$ and $K = 30$. For each set of parameters, 6 replicates were performed and averaged. For the GBD-backward model we varied two parameters, a and b , and fixed $d = 1$ and $c = 200$. The number of time units t was set to 5,000. We performed 15 replicates under each parameter combination in a grid of $(\frac{1}{a}, b)$ with $\frac{1}{a} \in [0.3, 3.5]$, every 0.2, and $b \in [0.01, 0.12]$, every 0.01.

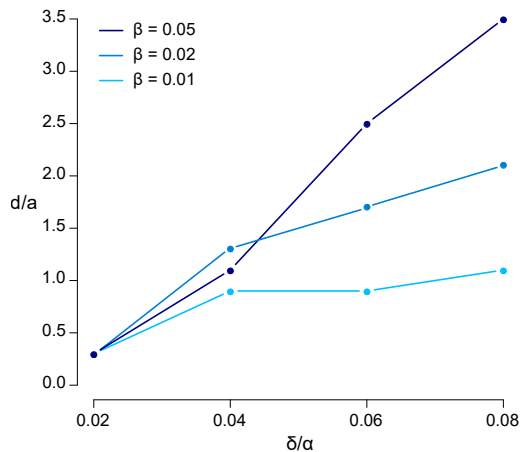
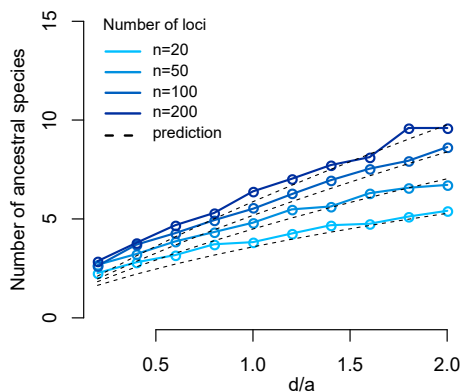


Figure 5: Evaluation of the GBD-backward model for a single sampled genome with n genes. Parameter settings: $a = 1$, $d \in [0.2, 2]$, every 0.2, and $n = 20, 50, 100$, and 200. The number of time units t was set to 10,000. We sampled the number of ancestral species every 500 time units starting at time $t = 5,000$, and averaged them for each simulation. For each set of parameters, 5 replicates were performed and averaged. A) Number of ancestral species depending on the number of genes n and on the ratio $\frac{d}{a}$, for one sampled genome. B) To assess the sampling consistency of our models, k lineages were randomly sampled. The number of ancestral species reported is the number of ancestral species of these k genes only.

A) Number of ancestral species



B) Sampling consistency (k=20)

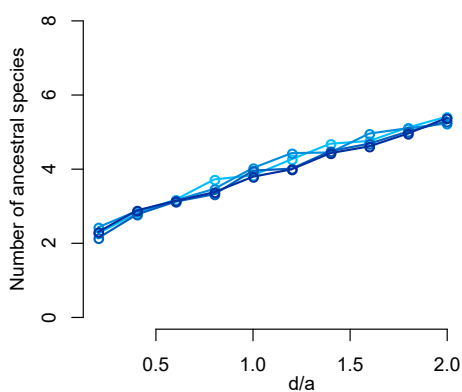


Figure 6: Billera-Holmes-Vogtmann (BHV) distances among sets of gene trees simulated under the gene-based diversification (GBD-backward) model. For each set of parameters, 15 simulations were performed (with $t = 5,000$, enough to reach the coalescence of all homologous genes) and the median BHV distances were calculated. A) Influence of the number of genes n (with $n = 5, 10$, and 20), of the number of species N (with $N = 6$ and 10), and of the ratio $\frac{d}{a}$ on the BHV distances. Parameter settings: $b = 0.05$, $d = 1$, $c = 200$, and $\frac{1}{a} = 0.3, 0.5, 0.9, 1.3, 1.7, 2.1, 2.5, 2.9, 3.3$. B) Influence of the *homologous attraction rate* b and of the ratio $\frac{d}{a}$ on the BHV distances. Parameter settings: $n = 10$, $N = 6$, $b = 0.01, 0.02, 0.05, 0.12$, $d = 1$, $c = 200$, and $\frac{1}{a} = 0.3, 0.5, 0.9, 1.3, 1.7, 2.1, 2.5, 2.9, 3.3$.

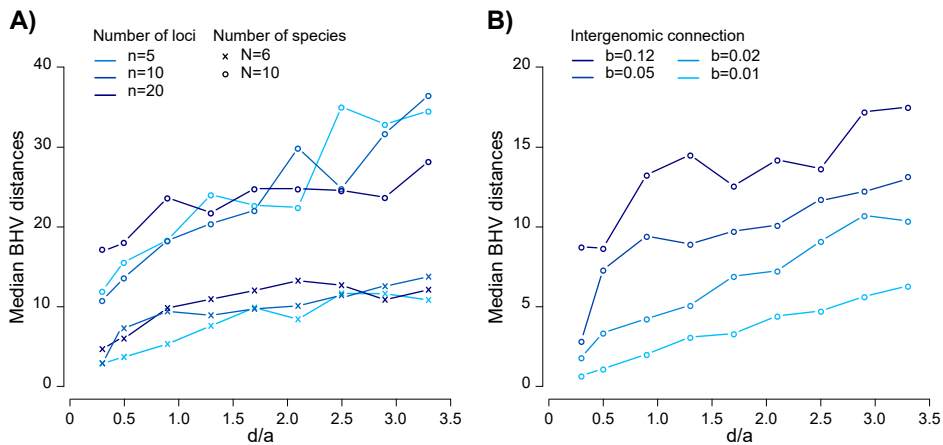


Figure 7: Bayesian phylogenetic reconstruction from simulated sequences under the GBD-backward model. We simulated 10 gene trees for 6 species under the GBD-backward model (with $\frac{b}{a} = 0.056$ and $\frac{d}{a} = 0.9$). The Bayesian phylogenetic analysis was performed with the program BEAST. The edges of the species tree (Bayesian analysis) are depicted by pipes in light gray. PP: posterior probabilities.

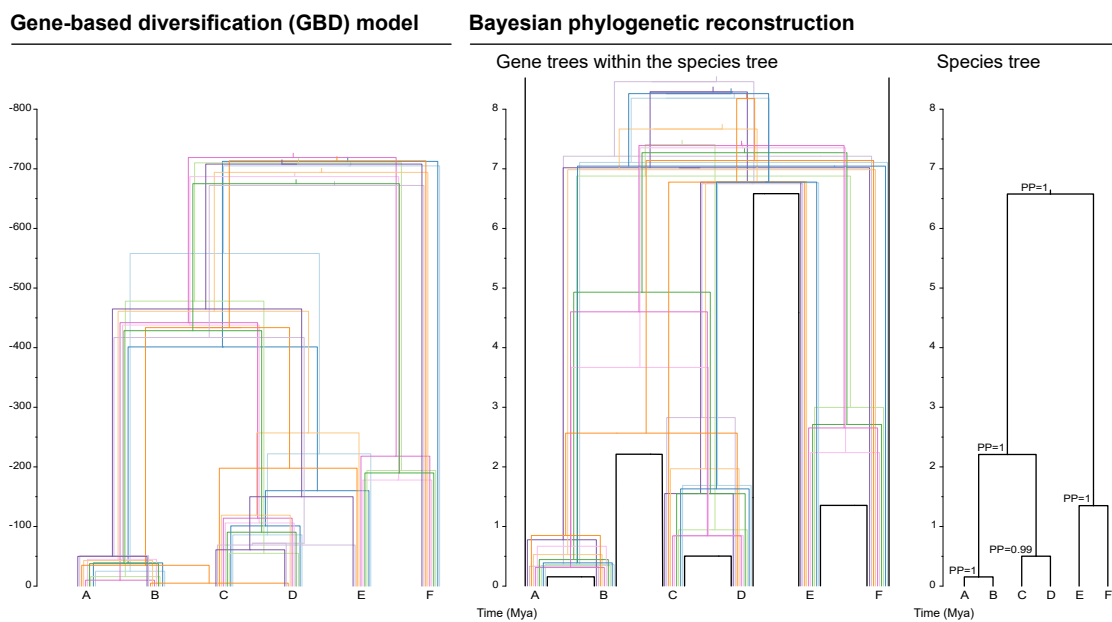
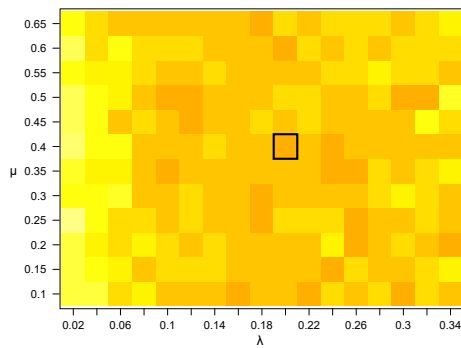


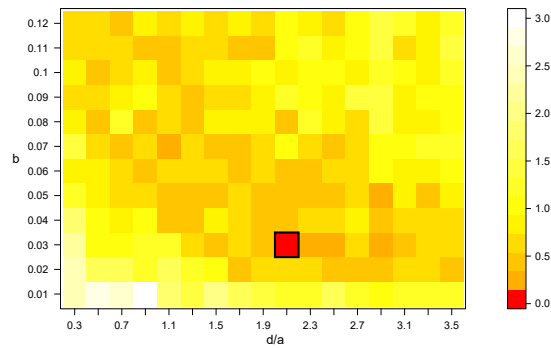
Figure 8: Minimization of the Kullback-Leibler (KL) divergence between empirical and simulated trees, *i.e.* between their distributions of BHV pairwise distances. Two parameters were optimized for each model. The *speciation* rate (λ) and the *extinction* rate (μ) for the multi-species coalescent (MSC) model (with coalescence rate set to 1). The *homologous attraction* b and the ratio of the *gene flow* rate over the *non-homologous attraction* rate ($\frac{d}{a}$) for the gene-based diversification (GBD-backward) model (with d set to 1). For each set of variables, 15 simulations were performed and averaged. The same color scale was used for each empirical data-set. For each optimization analysis, the cell for which we found the best fit between empirical and simulated trees (smallest KL divergence) is framed.

Bear data-set

Multi-species coalescent (MSC) model

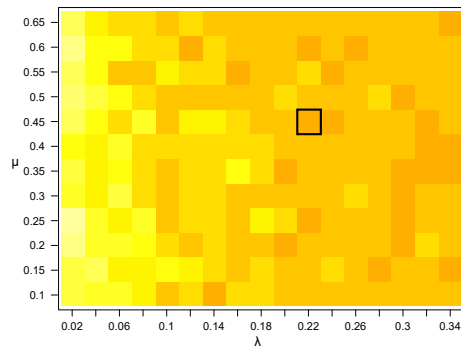


Gene-based diversification (GBD) model



Finch data-set

Multi-species coalescent (MSC) model



Gene-based diversification (GBD) model

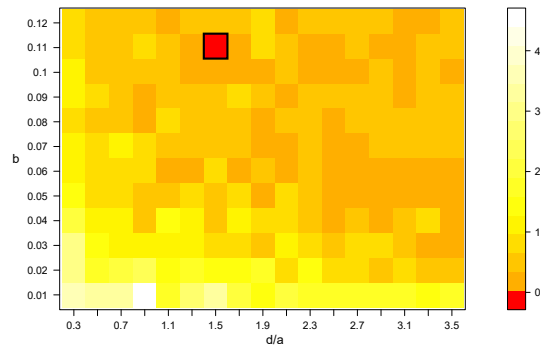
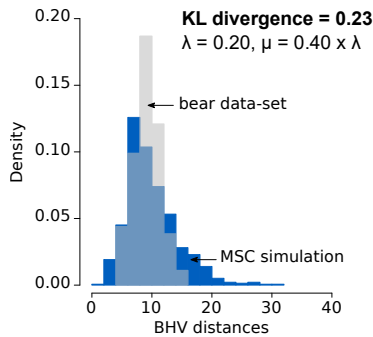


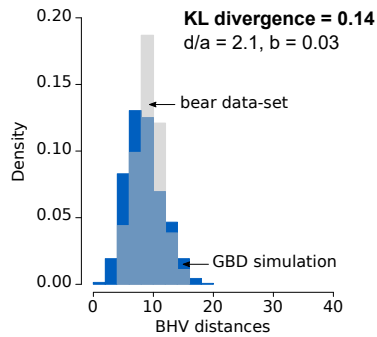
Figure 9: Best fit between empirical and simulated trees, *i.e.* between their distributions of BHV pairwise distances (selected cells of figure 8). For each set of variables, 15 simulations were performed and averaged. a : non-homologous attraction rate, b : homologous attraction rate, d : gene flow rate (set to 1), λ : speciation rate, μ : extinction rate, KL: Kullback-Leibler.

Bear data-set

Multi-species coalescent (MSC) model

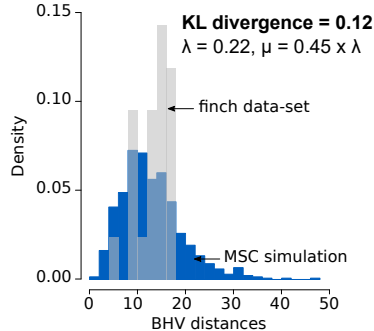


Gene-based diversification (GBD) model

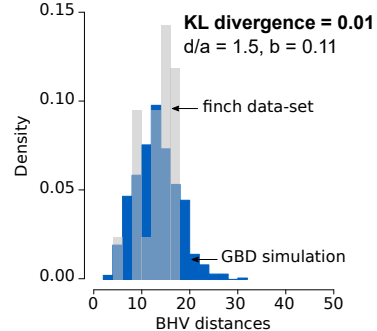


Finch data-set

Multi-species coalescent (MSC) model



Gene-based diversification (GBD) model



Supplementary figure 1

We tested our inference method, minimization of the difference (KL divergence) between the distributions of BHV distances,

on 8 simulated data sets (test data sets) with 10 replicates each. For each optimization analysis, the cell for which we found the best fit between the test trees and simulated trees (smallest KL divergence) is framed. The cross indicates the combination

parameters of the test data set.

