

1 Title: The genomic view of diversification

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

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

10
11 **ABSTRACT:** Evolutionary relationships between species are traditionally represented in the form of a
12 tree, called the species tree. The reconstruction of the species tree from molecular data is hindered
13 by frequent conflicts between gene genealogies. A standard way of dealing with this issue is to pos-
14 tulate the existence of a unique species tree where disagreements between gene trees are explained
15 by incomplete lineage sorting (ILS) due to random coalescences of gene lineages inside the edges of
16 the species tree. This paradigm, known as the multi-species coalescent (MSC), is constantly violated
17 by the ubiquitous presence of gene flow revealed by empirical studies, leading to topological incon-
18 gruenes of gene trees that cannot be explained by ILS alone. Here we argue that this paradigm
19 should be revised in favor of a vision acknowledging the importance of gene flow and where gene
20 histories shape the species tree rather than the opposite. We propose a new, plastic framework for
21 modeling the joint evolution of gene and species lineages relaxing the hierarchy between the species
22 tree and gene trees. As an illustration, we implement this framework in a mathematical model called
23 the genomic diversification (GD) model based on coalescent theory, with four parameters tuning repli-
24 cation, genetic differentiation, gene flow and reproductive isolation. We use it to evaluate the amount
25 of gene flow in two empirical data-sets. We find that in these data-sets, gene tree distributions are
26 better explained by the best fitting GD model than by the best fitting MSC model. This work should
27 pave the way for approaches of diversification using the richer signal contained in genomic evolution-
28 ary histories rather than in the mere species tree.

29
30 **Keywords:** coalescent theory, gene flow, gene tree, genomic diversification model, multi-species coa-
31 lescent, phylogeny, population genetics, speciation, species tree, reproductive isolation, introgression.

32 INTRODUCTION

33
34 The most widely used way of representing evolutionary relationships between contemporary species
35 is the so-called species tree, or phylogeny. The high efficiency of statistical methods using sequence
36 data to reconstruct species trees, hence called ‘molecular phylogenies’, led to precise dating of the

nodes of these phylogenies [34, 37, 82]. Notwithstanding the debatable accuracy of these datings, the use of time-calibrated phylogenies, sometimes called ‘timetrees’ [33], has progressively overtaken a view where phylogenies merely represent tree-like relationships between species in favor of a view where the timetree is the exact reflection of the diversification process [60, 67, 80]. In this view, the nodes of the phylogeny are consequently seen as punctual speciation events where one daughter species is instantaneously ‘born’ from a mother species. In this paper, we explore an alternative view of diversification, acknowledging that speciation is a long-term process [16, 42, 68] and not invoking any notion of mother-daughter relationship between species as done in the timetree view. This alternative view is gene-based rather than species-based, comparable with Wu’s genic view of speciation [85]. We use here the term ‘gene’ in the sense of “non-recombining locus”, *i.e.*, a region of the genome with a unique evolutionary history. Our view is meant in particular to accommodate the well-recognized existence of gene flow between incipient species, which persists during the speciation process and long after [50].

The timetree view of phylogenies does acknowledge that gene trees are not independent and may disagree with the species tree [47], but current methods jointly inferring gene trees and species tree rely on the following assumptions that we question in the next section: there is a unique species tree, the species tree shapes the gene trees and the species tree is the only factor mediating all dependences between gene trees (they are independent conditional on the species tree).

This view is materialized in a model called the ‘multispecies coalescent’ (MSC) [38] where conditional on the species tree, the evolutionary histories of genes follow independent coalescents constrained to take place within the hollow edges of the species tree. Many methods have been developed to estimate the species tree under the MSC, such as full likelihood methods (e.g. BEAST [34], BPP [88]) which average over gene trees and parameters [87], and the approximate or summary coalescent methods (e.g. ASTRAL [57], MP-EST [44], and STELLS [86]) which use a two-step approach: gene trees are first inferred and then combined to estimate the species tree that minimize conflicts among gene trees. Discordance between gene topologies is then explained, as a first approximation at least, by the intrinsic randomness of coalescences resulting in incomplete lineage sorting (ILS) (figure 1).

However, the presence of gene flow (hybridization, horizontal transfer) is now widely recognized between closely related species, and even between distantly related species [50]. Porous species boundaries, allowing for gene exchange because of incomplete reproductive isolation, are indeed regularly observed in diverse taxa such as amphibians [20, 65], arthropods [12], cichlids [84], cyprinids [6, 23, 24, 25, 79], insects [61, 64, 83], and even more frequently among bacteria [50, 78]. Long neglected, gene flow has recently been recognized as an important evolutionary driving force, through adaptive introgression or the formation of new hybrid taxa [1]. The ubiquity of genetic exchange across the Tree of Life between contemporary species suggests that gene flow has occurred many times in

the evolutionary past, and might actually be the most important cause of discrepancies between gene histories (e.g. [8, 11, 22, 36]) (figure 1). Accordingly, several extensions to the MSC model have been considered allowing for gene flow between species [39, 89]. These models acknowledge that species boundaries can be permeable at a few specific timepoints [32]. Unfortunately, because of the heavy computational cost of modeling the coalescent with gene flow, these methods are limited to small data-sets [89]. More importantly, they might not be appropriate to realistically model gene flow, given the frequency of gene flow across time and clades described in empirical studies [77]. Additionally, some of these methods, ASTRAL and MP-EST, might infer erroneous gene trees when gene flow is present [46]. These observations urge for novel approaches where gene flow is the rule rather than the exception.

To fill this void, we propose here an alternative model, that we call the genomic diversification (GD) model, framed with minimal assumptions arising from recent empirical evidence. Unlike the timetree view, our genomic view of diversification does not put the emphasis on the species tree (which in our model becomes a network rather than a tree) and assumes that gene trees shape the species tree (rather than the opposite).

THE GENOMIC VIEW OF DIVERSIFICATION

Gene flow and the questionable existence of a species genealogy

The biological species concept (BSC [53]) defines species as groups of interbreeding populations that are reproductively isolated from other groups. This definition postulates the non-permeability of species boundaries, which is contradicted by the growing body of evidence describing permeable or semi-permeable genomes, even between distantly related taxa. To integrate the possibility of gene flow into the definition of species, Wu [85] shifted the emphasis from isolation at the level of the whole genome to differential isolation at the gene level. Species are thus defined as differentially adapted groups for which inter-specific gene flow is allowed except for genes involved in differential adaptation (a well-defined form of divergence in which the alternative alleles have opposite fitness effects in the two groups) [85]. Because a fraction of the genome may still be exchanged after speciation is complete, a mosaic of gene genealogies is expected between divergent genomes [85]. Much evidence supports this prediction with the observation of highly conflicting gene trees, e.g. Darwin's finches [26, 28], sympatric sticklebacks [69, 73], Iberian barbels [24], and *Rhagoletis* species [2].

Accordingly, the notion of a species genealogy as the binary division of species into new independently evolving lineages in bifurcating phylogenetic trees, appears inappropriate. To avoid this misleading vision of speciation, we here wish to relax the species tree constraint by considering only gene genealogies as real genealogies, thereby laying aside, at least temporarily, the notion of species genealogy. To do so, we do not specify mother-daughter relationships between species, yet we postulate the existence of species at any time, and assume that we can unambiguously follow the

109 genealogies of genes (defined as non-recombining loci, as mentioned above).

110 Looking forward in time, genes belonging to two distinct individuals may find each other, in a next
111 generation, in the same genome because of recombination. The same process might occur with
112 two individuals belonging to different species under gene flow: this process, viewed in the backward
113 direction of natural time, is defined here as *disconnection* (figure 2).

114 A primary consequence of the presence of gene flow is to challenge the notion of a unique an-
115 cestral species. If all genes ancestral to species S have traveled through the same species in the
116 past, then species S has only one single ancestor species at any time. But because of gene flow (*i.e.*
117 *disconnection*), these genes may lie in different species living at a given time in the past, such that
118 species S can have several ancestral species at this time.

119

120 **Genome cohesion under continuous gene flow**

121 While some genes (e.g., genes involved in divergent adaptation) are hardly exchanged between
122 populations, other genes (e.g., neutral genes unlinked to genes under divergent selection) can be
123 subject to gene flow between different species [66, 85]. Gene flow can persist for long periods of
124 time, with evidence suggesting introgression events occurring on periods lasting up to 20 Myr [6, 24,
125 84]. Over time, genetic differences will accumulate in regions of low recombination and expand via
126 selective sweeps, leading eventually to complete reproductive isolation [85]. Accordingly, pairs of
127 species will likely exhibit greater genetic incompatibility through time, *i.e.* be less permeable to gene
128 flow, as has been observed for Iberian barbels [24], pea aphids [64], or salamanders [65]. In other
129 words, gene lineages remaining too long isolated within the same species decrease their ability to
130 introgress the genome of another species, a property that we name *genome cohesion* and which is
131 the consequence of spontaneous genetic differentiation.

132 Seen from the viewpoint of a present-day genome, genome cohesion means that ancestral lin-
133 eages of different genes in this genome cannot have spent too much time in the past in different
134 species. As time goes backward, this results in an *apparent* attractive force that brings together lin-
135 eages of genes belonging to the same present-day genome toward the same species, a phenomenon
136 we term *intragenomic connection*.

137 This has to be distinguished from the *coalescence* which refers to the point in time when the
138 lineages of *homologous genes* sampled from *different* genomes merge into a single lineage. To coa-
139 lesce, homologous genes must be located in the same individual, hence in the same species. We call
140 the *intergenomic connection* of two genomes sampled from *different* present-day species (figure 2)
141 the first event (in backward time) of migration of two homologous genes from each of these genomes
142 into the same species. Note that after coalescence (hence after intergenomic connection) of two
143 homologous lineages from the two different genomes, the resulting lineage is now common to these
144 two genomes. As a consequence of the mere intragenomic connection, going further back in time, all

other genes will then converge to the same species and further coalesce, until all homologous gene lineages have coalesced.

The genomic diversification (GD) model

We propose here a new plastic framework, derived from the genomic view of diversification described above, that acknowledges the importance of gene flow and relaxes the hierarchy between the species tree and gene trees. This model that we named the genomic diversification (GD) model, uses coalescent theory for modeling the joint evolution of gene and species lineages, reconciling phylogenomics with our current knowledge of species diversification. The GD model features only four parameters prescribing four processes affecting gene genealogies: *replication*, *genetic differentiation*, *introgression*, and *reproductive isolation*. In backward time, these four processes respectively become: *coalescence*, *intragenomic connection*, *disconnection* and *intergenomic connection*. To model the progressive isolation with ongoing gene flow, we assume only one event at a time (figure 2). This framework can be made more complex by letting the parameters depend on time, on the gene, or on any prescribed category of genes.

The GD model was implemented in R (<https://www.r-project.org>) and evaluated under different sets of parameters. We also applied it to two empirical multi-locus data-sets showing complex evolutionary patterns due to gene flow, each comprising six morphologically and ecologically distinct species, the Ursinae (a bear subfamily) [41] and the *Geospiza* clade (a genus of Darwin's finches) [19]. We estimated in particular 1) the relative amount of gene flow that has shaped each data-set, and 2) the corresponding average number of ancestral species.

MATERIAL AND METHODS

Parametrization of the GD model

At $t = 0$, n homologous genes are sampled in each of N sampled species. We will call a *block* at (backward) time t a (maximal) set of gene lineages that lie in the same species at time t , that are all ancestral to genes belonging to the same genome at $t = 0$. We now specify how we have parameterized the GD model. We follow the configuration of gene lineages into blocks in backward time, assuming a time-discrete Markov chain associated to the time-continuous chain with the following rates.

- **Intragenomic connection** (rate a). At any time t in the past, due to genome cohesion, each gene lineage at rate a independently, escapes from its block and chooses a target block with a probability proportional to the number of lineages (ancestral to genes belonging to the same genome at $t = 0$) harbored by the target block.
- **Intergenomic connection** (rate b). At any time t in the past, each pair of homologous genes (or preferably only of genes belonging to some previously prescribed category, like genes con-

181 contributing to reproductive isolation) find themselves in the same species at rate b independently.

182 • **Coalescence** (rate c). At any time t in the past, each pair of homologous genes lying within the
183 same species coalesces at rate c .

184 • **Disconnection** (rate d). At any time t in the past, as a result of gene flow, each gene lineage,
185 at rate d independently, escapes from its block and creates a new block (*i.e.*, gets into a species
186 harboring no other gene lineage) (figure 2). To model the introgression of bigger chunks of DNA,
187 we could alternatively assume that instead of one lineage, a given fraction of the lineages of a
188 block can simultaneously create a new block. We will not consider this possibility in the present
189 work.

190 We define the number of ancestral species of a given genome at time t , as the number of blocks
191 at time t containing the ancestral lineages to this genome. We considered a *time unit* to be equal to
192 the time elapsed between two events that we assumed to be constant for the sake of simplicity. In
193 this manuscript we wish to explore the impact of gene flow rather than ILS to explain gene tree con-
194 flicts, and thus scale a large c value (coalescence rate) so that coalescent events are instantaneous.
195 Therefore, only the parameters a , b , and d influence the gene genealogies.

196

197 **A single sampled genome**

198 We aimed to evaluate the variation in the number of ancestral species with gene flow. We per-
199 formed simulations for a single sampled genome containing n genes (with $n = 20, 50, 100, 200$), and
200 varied the relative amount of gene flow (*disconnection*) compared to genetic differentiation (*intra-*
201 *genomic connection*), ratio $\frac{d}{a}$ (with $a = 1$ and $d \in [0.2, 2]$, every 0.2). The number of time units t was
202 set to 10,000. We sampled the number of ancestral species every 500 time units starting at time
203 $t = 5,000$, and averaged these 11 values for each simulation. For each set of parameters, 5 replicates
204 were performed and averaged.

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

209

210 **A sample of several genomes**

211 When considering several sampled genomes of n genes, n gene genealogies are obtained for
212 a particular parameter setting. To characterize each set of gene genealogies, we employed a tree
213 comparison metric, the Billera-Holmes-Vogtmann (BHV) metric [3]. The BHV metric accounts for both
214 branch length and topological differences. This metric is a distance based on the concept of tree
215 space, a quadrant complex with quadrants sharing some faces. Two trees with the same topology

lie in the same quadrant, otherwise they lie in two distinct quadrants. At a common edge between two quadrants, the incongruent internal branches between trees have lengths equal to zero. Then a distance can be calculated between two rooted trees across these interconnected quadrants.

To compare trees that did not evolve on the same time scale, BHV distances were computed on re-scaled trees. For each set of gene trees issued from a single simulation or data-set, we rescaled all the trees so that the median of the most recent node depth is 1. We scaled the trees according to the median first coalescence among gene trees because in our model, the first coalescence initiates the intragenomic connection between genomes of different species, and hence coalescence of all the remaining homologous genes.

We evaluated the influence of the number of genes n (with $n = 5, 10, 20$), of the number of species N (with $N = 6, 10$), and of the relative amount of gene flow $\frac{d}{a}$ (with $d=1$ and $\frac{1}{a} = 0.3, 0.5, 0.9, 1.3, 1.7, 2.1, 2.5, 2.9, 3.3$) on gene tree diversity (BHV distances) (figure 4A). The other parameters were fixed, with $b = 0.05$ and $c = 200$.

For the same values of $\frac{d}{a}$ and c , and for $n = 10$, $N = 6$, we also evaluated the influence of the *intergenomic connection* rate b (with $b = 0.01, 0.02, 0.05, 0.12$) on gene tree diversity (BHV distances) (figure 4B).

The GD model versus the MSC model

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

A set of 10 gene trees was simulated under the GD model (with $N = 6$, $b = 0.05$, $\frac{1}{a} = 0.9$, and $d = 1$) (figure 5). We simulated DNA sequences (package 'PhyloSim' in R [75]) corresponding to each of the 10 gene trees with model of DNA evolution estimated by modeltest (function 'modelTest', package 'phangorn' in R [72]) for the TRAPPC10 intron of the bear data-set detailed below [41]: HKY model, rate matrix: $a = 1.00$, $b = 5.29$, $c = 1.00$, $d = 1.00$, $e = 5.29$, $f = 1.00$, base frequencies: 0.26, 0.19, 0.21, 0.34. Prior to simulating the sequences, the 10 gene trees were scaled to the TRAPP10 intron phylogenetic tree length (built with RaXML 8.1.11 [81] assuming GTR (general time reversible) model with 1,000 bootstrap replicates).

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

252

253 Inferences from empirical data-sets

254 Empirical data-sets

255 The amount of gene flow that has shaped the two empirical data-sets was estimated by comparing
 256 the distributions of their pairwise gene tree distances with those of simulated trees. The first data-set
 257 comprised 14 autosomal introns for 6 bear species (*Helarctos malayanus*, *Melursus ursinus*, *Ursus*
 258 *americanus*, *U. arctos*, *U. maritimus*, and *U. thibetanus*) and 2 outgroups (*Ailuropoda melanoleuca*
 259 and *Tremarctos ornatus*) [41]. The sequences were downloaded from GenBank (supplementary table
 260 S1). As in Kutschera et al. [41], all variation within and among individuals was collapsed into one
 261 single 50% majority-rule-consensus sequence for each of the 8 species. The phylogenetic trees were
 262 built with the program BEAST v. 1.8.3. [13], with the parameters used by the authors of [41]: Yule prior
 263 to model the branching process, strict clock, a normal prior on substitution rates (0.001 ± 0.001) (mean
 264 \pm SD), minimum age of 11.6 My for the divergence of *A. melanoleuca* from other bears (exponential
 265 prior: mean= 0.5; offset= 11.6), and 10 million generations with sampling every 1000 generations. The
 266 models of DNA evolution were estimated by modeltest (function 'modelTest', package 'phangorn' in R
 267 [72]) (supplementary table S2). The monophyly of the ingroup and the topology among the outgroups
 268 were constrained according to the topology depicted in Kutschera et al. [41].

269 The second data-set comprised 7 nuclear markers for 6 finch species (*Geospiza conirostris*,
 270 *G. fortis*, *G. fuliginosa*, *G. magnirostris*, *G. scandens*, and *G. septentrionalis*) and 2 outgroups (*Ca-*
 271 *marhynchus psittacula* and *Platyspiza crassirostris*) [19]. The sequences were downloaded from
 272 GenBank (supplementary table S3). The phylogenetic trees were built with the program BEAST v.
 273 1.8.3. [13] with the parameters used by Farrington et al. [19]: coalescent constant size prior to model
 274 the branching process, strict clock, substitution rate equal to 1, specific models of DNA evolution de-
 275 fined by the authors (supplementary table S2), and 10 million generations with sampling every 1000
 276 generations. The monophyly of the ingroup and the topology among the outgroups were constrained
 277 according to the topology depicted in [19].

278

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

280 We optimized the MSC model for $N = 6$ species by varying two parameters, the speciation rate
 281 λ and the extinction rate μ , and fixing the coalescence rate to 1. Birth-death trees of 6 tips (function
 282 'sim.bdtree', package 'geiger' in R) were simulated in a grid of $(\lambda, \mu = m\lambda)$ with $\lambda \in [0.02, 0.34]$, every
 283 0.02, and $m \in [0.1, 0.65]$, every 0.05. Because we simulated small trees (6 tips), the degree of variation
 284 between trees simulated with the same parameters was high. Therefore for each value of (λ, μ) we
 285 randomly selected 15 species trees for which the crown age did not differ by more than 2.5% from the
 286 expected crown age. Next, we simulated 10 gene genealogies for each species tree (coalescence
 287 rate fixed to 1).

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

Estimation of parameters under the genomic diversification (GD) model

Equivalently, we optimized the GD model for $N = 6$ by varying two parameters, here a and b , and fixing $d = 1$ and $c = 200$ (recall c is given a sufficiently large value that coalescences are instantaneous). Since increasing n has no effect on BHV distances (see above and figure 4), we simulated genomes with $n = 10$ genes. The number of time units t was set to 5,000, which guarantees the coalescence of all homologous genes. 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.

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

RESULTS

A single sampled genome

Let us consider the case of $N = 1$ sampled genome containing n genes. We let $A(t) = (A_1(t), \dots, A_n(t))$ denote the sorting of genes into ancestral species t units of time before the present. More precisely, $A_k(t)$ denotes the number of ancestral species containing k gene lineages, so that $n = \sum_{k=1}^n k A_k(t)$ and $S(t) = \sum_{k=1}^n A_k(t)$ is the total number of species at t ancestral to the sampled genome. For each $\varepsilon \in (0, 1]$, 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 at least a fraction ε of the genome (with $\lceil x \rceil$ denoting the smallest integer larger than x). All stationary quantities will be denoted by the same symbols, replacing t with ∞ .

The transition rates can be specified as follows in terms of the configuration of gene lineages into blocks (*i.e.*, ancestral species). For each pair of blocks containing (j, k) lineages, intragenomic connection occurs at rate ajk and results in the configuration $(j - 1, k + 1)$. For each block containing j lineages, disconnection occurs at rate dj and results in the block losing one lineage; simultaneously a new block containing 1 single lineage is created. These are exactly the same rates as in the well-known Moran model with mutation under the infinite-allele model [58], replacing 'block' with 'allele', 'connection' by 'resampling' (simultaneous birth from one of the j carriers of a given allele and death of one of the k carriers of another given allele) and 'disconnection' with 'mutation' (mutation appearing in one of the j carriers of a given allele into a new allele never existing before). For this Moran model,

- the total population size is n ;
- at rate a for each oriented pair of individuals independently, the first individual of the pair gives birth to a copy of herself and the second individual of the pair is simultaneously killed;
- 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 [14, 17, 18]. Expectations of this distribution are:

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

so that

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

and

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

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

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

330

At stationarity, and particularly for large values of $\frac{d}{a}$, the mean number of ancestral species $S(\infty)$ obtained from simulations was equal to the mathematical prediction (figure 3A). In particular, the mean number of ancestral species at stationarity increases with $\frac{d}{a}$.

An additional key feature of this model is *sampling consistency*. In words, the history of a sample of k genes taken from a genome of n genes does not depend on n . This property can again be deduced from the representation of our model in terms of the better known Moran model. Indeed, the dynamics of a sample of k individuals in the Moran model does not depend on the population size, as can be seen from the so-called lookdown construction [15]. The simulations performed with k genes randomly sampled from each genome of n genes, are in agreement with this claim of sampling consistency: the number of ancestral species at stationarity $\mathbb{E}(S(\infty))$ is independent of the number of genes n (figure 3B).

342

343 A sample of several genomes

Using simulations, we evaluated the GD model for several sampled genomes ($N > 1$) under several combinations of parameters. As expected gene tree diversity, measured by BHV distances,

345

increased with $\frac{d}{a}$, *i.e.* the relative amount of gene flow, and with the number of species N . Conversely our results showed that the number of genes n had no effect on distances (figure 4A). This last result, the lack of influence of n on gene tree diversity, is of particular interest, because one usually has only access to a fraction of a genome. It shows that regardless of the number of genes sampled, the resulting gene tree diversity will remain the same as long as gene trees have been shaped by processes with similar parameter values.

Our results also showed that as the *intergenomic connection* rate b decreases, and for the same $\frac{d}{a}$, gene trees were more similar (lower BHV distances) (figure 4B). When a long period of time elapses between two intergenomic connection events (low b), all the genes belonging to the two genomes that have started to coalesce, have enough time to converge toward the same species, and thus coalesce before the next intergenomic connection event, in spite of gene flow.

GD versus MSC: ignoring gene flow may lead to mistaken phylogenetic inferences

When evaluating the ability of MSC model to deal with gene flow, we found a strong support (posterior probabilities > 0.90) for all the nodes of the Bayesian species tree even if the individual gene trees of the GD model did not corroborate this topology (figure 5). For example, 7 out of 10 gene trees modeled under the GD model support the connection between the species E and the species C and D, and only 3 the direct relationship between the species E and F. Whereas the Bayesian tree strongly supports the clade (E,F) with a posterior probability equal to 1, and considers all the connections between E and (C,D) to be due to ancestral polymorphism (*i.e.*, ILS). Moreover because gene trees are constrained in the species tree (MSC model), the coalescences between genes of E and (C,D) must take place after the species tree coalescence, therefore these coalescences are timed around 7 My instead of 2 My according to the GD tree. Failing to recognize that that gene flow may have shaped gene genealogies, hence DNA sequences, can result in important topological and dating errors.

Inferences from empirical data-sets support the GD model

To find the best set of parameters, we minimized the Kullback-Leibler (KL) divergence between the distributions of BHV pairwise distances of empirical and simulated trees (figure 6). Under the multi-species coalescent (MSC) model, the most likely set of parameters was $\mu = 0.4 \times \lambda$ and $\lambda = 0.2$ (KL divergence = 0.23) for the bears, and $\mu = 0.45 \times \lambda$ and $\lambda = 0.22$ for the finches (KL divergence = 0.12). We noted longer tailed distributions for the distances between trees modeled under the MSC model than for the empirical data-sets (figure 7). This skewed distribution obtained with the MSC model explains why we did not detect a sharp peak in the optimization landscape for the MSC model (figure 6).

Under the genomic diversification (GD) model, the most likely set of parameters 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 the finches (KL

divergence = 0.01) (figure 6). Contrary to the MSC model, the distributions of the distances between trees modeled under the GD model or empirical trees did not show, or to a lesser degree, a long tail (figure 7), explaining why we could detect a sharp peak in the optimization landscape for the MSC model (figure 6).

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

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

DISCUSSION

Within species, gene flow allows the maintenance of species cohesion in the face of genetic differentiation [59, 76], preventing genetic isolation of populations and the subsequent emergence of reproductive barriers leading to speciation [10]. Among species, the existence of gene flow challenges the notion of a species genealogy as well as the current concepts of species. Indeed, if gene flow is as pervasive as recent empirical studies suggest [8, 11, 22, 36], the genealogical history of species should be represented as a phylogenetic network encompassing the mosaic of gene genealogies. Similarly, it seems very conservative to delineate species based on the widely used biological species concept (reproductive isolation) [53], or phylogenetic species concept (reciprocal monophyly) [62]. Because of the ubiquity of gene flow, which can persist for several millions of years after the lineages have started to diverge (*i.e.*, onset of speciation) [4, 48], species should be rather defined by their capacity to coexist without fusion in spite of gene flow [49, 70].

The simplified view of diversification, consisting in representing lineages splitting instantaneously into divergent lineages with no interaction (gene exchange) after the split, has been preventing evolutionary biologists from fully apprehending diversification at the genomic level and from correctly interpreting discrepancies between gene histories. Indeed, conflicting gene trees make the interpretation of their evolutionary history difficult. However, we argue that phylogenetic incongruence among gene trees should not be considered as a nuisance, but rather as a meaningful biological signal revealing some features of the dynamics of genetic differentiation and of gene flow through time and across clades. Current phylogenetic methods rely on the assumption that gene trees are constrained within the species tree, and that gene flow occurs infrequently between species. For many data-sets such as sequence alignments of genomes sampled from young clades, such methods could lead

to an evolutionary misinterpretation of gene trees, and in the worst case to species trees with high node support while the gene trees had very different evolutionary histories (see figure 5). These observations urge for a change of paradigm, where gene flow is fully part of the diversification model. To consider the ubiquity of gene flow across the Tree of Life described by many recent studies, we have developed a new framework focusing on gene genealogies and relaxing the constraints inherent to the MSC paradigm. This framework is materialized in a mathematical model that we named the genomic diversification (GD) model.

The GD model

Under the GD model, gene genealogies are governed by four parameters corresponding to four biological processes, *coalescence* (replication), *intragenomic connection* (genetic differentiation), *intergenomic connection* (reproductive isolation), and *disconnection* (introgression) (figure 2).

Intergenomic connection corresponds to finding the most recent common ancestor of the two species at the genomic level. The time spent between intergenomic connections depends crucially on the (phylogenetic distance of the) species sampled at the present. *Disconnection* corresponds to the introgression of genetic material from one species into another species, which rate scales with the intensity of gene flow. *Intragenomic connection* models genetic differentiation. The slower genes accumulate mutations and differentiate, the more time can be spent by gene lineages in different species. Hence when genomes differentiate slowly, the rate of intragenomic connection is low.

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

Higher amount of gene flow (*disconnection*) and reduced time to untangle gene genealogies before the connection of two other genomes (*intergenomic connection*) increase the diversity among trees. Conversely, when homologous genes coalesce faster (*coalescence*) and genes converge faster toward the species harboring the other genes of their genome (*intragenomic connection*) a lower diversity among trees is expected.

After evaluating this model under various sets of parameters, we applied it to analyze two empirical multi-locus data-sets for which gene tree conflicts have obscured the evolutionary history.

Gene flow among bears and among finches

Our results showed support for the hypothesis that gene flow has shaped the gene trees of bears and finches (figure 7). For the bear data-set, we found that each species had on average in the past about 4.8 ancestral species carrying at least 10% of its present genome (equation (2)). This result is in line with previous studies reporting gene flow between pairs of bear species [7, 31, 41, 45, 55].

Moreover, a recent phylogenomic study (869 Mb divided into 18,621 genome fragments) confirmed the existence of gene flow between sister species as well as between more phylogenetically distant species [40]. They used the D -statistics (gene flow between sister species) and D_{FOIL} -statistics (gene flow among ancestral lineages [63]) to detect gene flow among the 6 bear species. Using their results, for each pair of species ij among the N species, we determined if the species j has contributed ($g_{ij} = 1$) or not ($g_{ij} = 0$) to the genome of the species i (with $g_{ii} = 1$), and calculated the average number of ancestral species S as follow:

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

We found on average 5.3 ancestral species for each of the Ursinae bears [40], close to the estimate obtained with the GD model (4.8).

We detected lower gene flow among finches than among bears. Each finch species had on average in the past 3.4 ancestral species (for the subsample of gene trees analyzed here), which is also consistent with the extensive evidence that many species hybridize on several islands [21, 27, 29, 30, 71]. Because of gene flow very little genetic structure was detected by a Bayesian population structure analysis, only 3 genetic populations among the 6 *Geospiza* species [19]. Each of the 2 species, *G. magnirostris* and *G. scandens*, were mostly characterized by a single genetic population, therefore had about 1 ancestral species each. Conversely 4 *Geospiza* species shared the same genetic population, suggesting 4 ancestral species for each of these 4 species. Taking together these results roughly indicate that each of the 6 *Geospiza* species had in average 3 ancestral species, in line with the GD estimate (3.4).

We showed here that strictly bifurcating lineage-based models do not adequately capture complex evolutionary patterns at the species level. On the contrary, a model relaxing species boundaries and accounting for gene flow, like the GD model, better reproduced the complex history of gene genealogies under continuous gene flow. Note that we considered a simple scenario with no ILS and statistically exchangeable genes resulting in a model with only three parameters, but given the simplicity and the flexibility of our model, many extensions may be considered to address scenarios that could not have been considered previously, opening up new perspectives in the study of speciation and macro-evolution.

Gene flow: an evolutionary force driving diversification

Species diversification requires genetic variation among organisms, introduced by mutations and structural variation, upon which natural selection and drift can act by influencing the sorting of offspring and the survival of organisms [70]. Recently, gene flow has also been mentioned as another potential source of genetic variation [52], and more particularly in the case of adaptive radiations [9, 43, 54, 74]. Hybrid zones act as filters, preventing the introgression of deleterious genes while

allowing advantageous or neutral genes to cross the species boundaries [52]. Newly acquired genes will then be a source of variation [52], by providing evolutionary adaptive shortcuts (beneficial genes) or greater adaptability once in the genetic pool of the introgressed species (neutral markers) [52]. The introgressed species then has a wider range of potentially adaptive allelic variants, allowing it to diversify rapidly if the opportunity arises. Accordingly important gene flow should be detected prior to an adaptive radiation. This hypothesis is supported by empirical evidence, but has only been tested under limited conditions [9, 43, 54, 74]. The model proposed here constitutes a great opportunity to investigate more systematically how gene flow is distributed throughout the phylogenies and how it can influence the frequency of adaptive radiations.

Evolutionary dynamics along the genome

Along the genome, gene flow is not expected to be uniformly distributed either. Incongruent gene trees should reveal genes that have evolved more slowly. Indeed, because of the genome cohesion force, genes evolving slower will be able to stay longer in different species. Conversely, congruent gene trees should reveal genomic regions not subject to gene flow, as genomic regions under strong selective differentiation [32, 35]. This framework could thus be used to evaluate how gene flow varies along the genome and to explore the genomic architecture of species barriers. Indeed some regions, as sexual chromosomes or low recombination genomic regions, are expected to be more differentiated and hence to undergo less gene flow (e.g. *Heliconius* species [51]). In order to distinguish between genes and to reduce potential errors in parameter estimation, data may be grouped by gene class (statistical binning) using a method aiming to evaluate whether two genes are likely to have the same tree (linked sites) or the same tree in distribution (statistical exchangeability) [56].

Perspectives

Models and methods inferring macro-evolutionary history from phylogenetic trees, such as speciation and extinction rates, trait evolution, and ancestral character reconstruction, have become increasingly complex [60, 67, 80]. Yet, the raw material used by these methods is often reduced to the species tree, which can be viewed as a summary statistic of the information contained in the genome. We argue here that a valuable amount of additional signal, not accessible in phylogenetic trees, is contained in gene trees, and is directly informative about the diversification process. Indeed, because genetic differentiation and gene flow impact each gene differently, genes may have experienced very different evolutionary trajectories.

In order to make use of the entire information conveyed by gene trees, we propose here a new approach to tackle the diversification process, the genomic view of diversification, under which gene trees shape the species tree rather than the opposite. This approach aims at better depicting the intricate evolutionary history of species and genomes. We hope that this view of diversification will

pave the way for future developments in the perspective of inferring diversification processes directly from genomes rather than from their summary into one single species tree. One of the challenges in this direction will be to propose finer inference methods than the crude one used here, based on a single summary statistic, the BHV distances.

SUPPLEMENTARY MATERIAL

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

ACKNOWLEDGMENTS

The authors thank the *Center for Interdisciplinary Research in Biology* (Collège de France, CNRS) for funding. JM is funded by LabEx MemoLife, project *Genomics of Diversification*.

References

- [1] ABBOTT, R., ALBACH, D., ANSELL, S., ARNTZEN, J. W., BAIRD, S. J. E., BIERNE, N., BOUGHMAN, J., BRELSFORD, A., BUERKLE, C. A., BUGGS, R., BUTLIN, R. K., DIECKMANN, U., EROUKHMANOFF, F., GRILL, A., CAHAN, S. H., HERMANSEN, J. S., HEWITT, G., HUDSON, A. G., JIGGINS, C., JONES, J., KELLER, B., MARCZEWSKI, T., MALLET, J., MARTINEZ-RODRIGUEZ, P., MÖST, M., MULLEN, S., NICHOLS, R., NOLTE, A. W., PARISOD, C., PFENNIG, K., RICE, A. M., RITCHIE, M. G., SEIFERT, B., SMADJA, C. M., STELKENS, R., SZYMURA, J. M., VÄINÖLÄ, R., WOLF, J. B. W., AND ZINNER, D. Hybridization and speciation. *Journal of Evolutionary Biology* 26, 2 (2013), 229–246.
- [2] BERLOCHER, S. H. Radiation and divergence in the *Rhagoletis pomonella* species group: inferences from allozymes. *Evolution* 54, 2 (2000), 543–557.
- [3] BILLERA, L. J., HOLMES, S. P., AND VOGTMANN, K. Geometry of the space of phylogenetic trees. *Advances in Applied Mathematics* 27, 4 (2001), 733–767.
- [4] BOLNICK, D. I., NEAR, T. J., AND NOOR, M. Tempo of hybrid inviability in centrarchid fishes (teleostei: centrarchidae). *Evolution* 59, 8 (2005), 1754–1767.
- [5] BOUCKAERT, R., HELED, J., KÜHNERT, D., VAUGHAN, T., WU, C.-H., XIE, D., SUCHARD, M. A., RAMBAUT, A., AND DRUMMOND, A. J. BEAST 2: A Software Platform for Bayesian Evolutionary Analysis. *PLOS Computational Biology* 10, 4 (2014), e1003537.
- [6] BUONERBA, L., ZACCARA, S., DELMASTRO, G. B., LORENZONI, M., SALZBURGER, W., AND GANTE, H. F. Intrinsic and extrinsic factors act at different spatial and temporal scales to shape

- 556 population structure, distribution and speciation in Italian *Barbus* (Osteichthyes: Cyprinidae).
557 *Molecular Phylogenetics and Evolution* 89 (2015), 115–129.
- 558 [7] CAHILL, J. A., GREEN, R. E., FULTON, T. L., STILLER, M., JAY, F., OVSYANIKOV, N.,
559 SALAMZADE, R., JOHN, J. S., STIRLING, I., AND SLATKIN, M. Genomic evidence for island
560 population conversion resolves conflicting theories of polar bear evolution. *PLoS genetics* 9, 3
561 (2013), e1003345.
- 562 [8] CLARK, A. G., AND MESSER, P. W. Conundrum of jumbled mosquito genomes. *Science* 347,
563 6217 (2015), 27–28.
- 564 [9] CONSORTIUM, H. G. Butterfly genome reveals promiscuous exchange of mimicry adaptations
565 among species. *Nature* 487, 7405 (2012), 94–98.
- 566 [10] COYNE, J. A., AND ORR, H. A. *Sympatric speciation. Speciation*. Sinauer Associates, Sunder-
567 land, MA, 2004.
- 568 [11] CUI, R., SCHUMER, M., KRUESI, K., WALTER, R., ANDOLFATTO, P., AND ROSENTHAL, G. G.
569 Phylogenomics reveals extensive reticulate evolution in *Xiphophorus* fishes. *Evolution* 67, 8
570 (2013), 2166–2179.
- 571 [12] DE BUSSCHERE, C., HENDRICKX, F., VAN BELLEGHEM, S. M., BACKELJAU, T., LENS, L., AND
572 BAERT, L. Parallel habitat specialization within the wolf spider genus Hogna from the Galápagos.
573 *Molecular ecology* 19, 18 (2010), 4029–4045.
- 574 [13] DRUMMOND, A. J., SUCHARD, M. A., XIE, D., AND RAMBAUT, A. Bayesian phylogenetics with
575 BEAUti and the BEAST 1.7. *Molecular biology and evolution* 29, 8 (2012), 1969–1973.
- 576 [14] DURRETT, R. *Probability Models for DNA Sequence Evolution*. Springer, 2008. Google-Books-
577 ID: o4_bMHY7jFoC.
- 578 [15] ETHERIDGE, A. *Some Mathematical Models from Population Genetics: École D'Été de Probabil-
579 ités de Saint-Flour XXXIX-2009*. Springer Science & Business Media, 2011. Google-Books-ID:
580 mil9tdPCFdUC.
- 581 [16] ETIENNE, R. S., MORLON, H., AND LAMBERT, A. Estimating the duration of speciation from
582 phylogenies. *Evolution* 68, 8 (2014), 2430–2440.
- 583 [17] EWENS, W. J. The sampling theory of selectively neutral alleles. *Theoretical Population Biology*
584 3, 1 (1972), 87–112.
- 585 [18] EWENS, W. J., AND TAVARÉ, S. Ewens Sampling Formula. In *Encyclopedia of Statistical Sci-
586 ences*. American Cancer Society, 2006.

- [19] FARRINGTON, H. L., LAWSON, L. P., CLARK, C. M., AND PETREN, K. The evolutionary history of Darwin's finches: speciation, gene flow, and introgression in a fragmented landscape. *Evolution* 68, 10 (2014), 2932–2944.
- [20] FONTENOT, B. E., MAKOWSKY, R., AND CHIPPIINDALE, P. T. Nuclear–mitochondrial discordance and gene flow in a recent radiation of toads. *Molecular Phylogenetics and Evolution* 59, 1 (2011), 66–80.
- [21] FREELAND, J. R., AND BOAG, P. T. The mitochondrial and nuclear genetic homogeneity of the phenotypically diverse darwin's ground finches. *Evolution; International Journal of Organic Evolution* 53, 5 (1999), 1553–1563.
- [22] GALLUS, S., JANKE, A., KUMAR, V., AND NILSSON, M. A. Disentangling the relationship of the Australian marsupial orders using retrotransposon and evolutionary network analyses. *Genome biology and evolution* 7, 4 (2015), 985–992.
- [23] GANTE, H. F., COLLARES-PEREIRA, M. J., AND COELHO, M. M. Introgressive hybridisation between two Iberian *Chondrostoma* species (Teleostei, Cyprinidae) revisited: new evidence from morphology, mitochondrial DNA, allozymes and NOR-phenotypes. *Folia Zoologica* 53, 4 (2004), 423.
- [24] GANTE, H. F., DOADRIO, I., ALVES, M. J., AND DOWLING, T. E. Semi-permeable species boundaries in Iberian barbels (*Barbus* and *Luciobarbus*, Cyprinidae). *BMC evolutionary biology* 15, 1 (2015), 111.
- [25] GANTE, H. F., SANTOS, C. D., AND ALVES, M. J. Phylogenetic relationships of the newly described species *Chondrostoma olisiponensis* (Teleostei: Cyprinidae). *Journal of Fish Biology* 76, 4 (2010), 965–974.
- [26] GRANT, B. R., AND GRANT, P. R. Hybridization and speciation in Darwin's finches: the role of sexual imprinting on a culturally transmitted trait. *Endless forms: species and speciation* (1998), 404–422.
- [27] GRANT, P. R., AND GRANT, B. R. Phenotypic and genetic effects of hybridization in darwin's finches. *Evolution; International Journal of Organic Evolution* 48, 2 (1994), 297–316.
- [28] GRANT, P. R., AND GRANT, B. R. Speciation and hybridization in island birds. *Phil. Trans. R. Soc. Lond. B* 351, 1341 (1996), 765–772.
- [29] GRANT, P. R., AND GRANT, B. R. Hybridization, Sexual Imprinting, and Mate Choice. *The American Naturalist* 149, 1 (1997), 1–28.

- [30] GRANT, P. R., GRANT, B. R., AND PETREN, K. Hybridization in the recent past. *The American Naturalist* 166, 1 (2005), 56–67.
- [31] HAILER, F., KUTSCHERA, V. E., HALLSTRÖM, B. M., KLASSERT, D., FAIN, S. R., LEONARD, J. A., ARNASON, U., AND JANKE, A. Nuclear genomic sequences reveal that polar bears are an old and distinct bear lineage. *Science* 336, 6079 (2012), 344–347.
- [32] HARRISON, R. G., AND LARSON, E. L. Hybridization, Introgression, and the Nature of Species Boundaries. *Journal of Heredity* 105, S1 (2014), 795–809.
- [33] HEDGES, S. B., AND KUMAR, S. *The Timetree of Life*. OUP Oxford, 2009. Google-Books-ID: 9rt1c1hl49MC.
- [34] HELED, J., AND DRUMMOND, A. J. Bayesian Inference of Species Trees from Multilocus Data. *Molecular Biology and Evolution* 27, 3 (2010), 570–580.
- [35] JANOUŠEK, V., MUNCLINGER, P., WANG, L., TEETER, K. C., AND TUCKER, P. K. Functional organization of the genome may shape the species boundary in the house mouse. *Molecular biology and evolution* 32, 5 (2015), 1208–1220.
- [36] JÓNSSON, H., SCHUBERT, M., SEGUIN-ORLANDO, A., GINOLHAC, A., PETERSEN, L., FUMAGALLI, M., ALBRECHTSEN, A., PETERSEN, B., KORNELIUSSEN, T. S., VILSTRUP, J. T., LEAR, T., MYKA, J. L., LUNDQUIST, J., MILLER, D. C., ALFARHAN, A. H., ALQURAISHI, S. A., ALRASHEID, K. A. S., STAGEGAARD, J., STRAUSS, G., BERTELSEN, M. F., SICHERITZ-PONTEN, T., ANTČZAK, D. F., BAILEY, E., NIELSEN, R., WILLERSLEV, E., AND ORLANDO, L. Speciation with gene flow in equids despite extensive chromosomal plasticity. *Proceedings of the National Academy of Sciences* 111, 52 (2014), 18655–18660.
- [37] KISHINO, H., THORNE, J. L., AND BRUNO, W. J. Performance of a divergence time estimation method under a probabilistic model of rate evolution. *Molecular Biology and Evolution* 18, 3 (2001), 352–361.
- [38] KNOWLES, L. L., AND KUBATKO, L. S. *Estimating species trees: practical and theoretical aspects*. John Wiley and Sons, 2011.
- [39] KUBATKO, L. S. Identifying hybridization events in the presence of coalescence via model selection. *Systematic Biology* 58, 5 (2009), 478–488.
- [40] KUMAR, V., LAMMERS, F., BIDON, T., PFENNINGER, M., KOLTER, L., NILSSON, M. A., AND JANKE, A. The evolutionary history of bears is characterized by gene flow across species. *Scientific Reports* 7 (2017), 46487.

- [41] KUTSCHERA, V. E., BIDON, T., HAILER, F., RODI, J. L., FAIN, S. R., AND JANKE, A. Bears in a forest of gene trees: phylogenetic inference is complicated by incomplete lineage sorting and gene flow. *Molecular Biology and Evolution* 31, 8 (2014), 2004–2017.
- [42] LAMBERT, A., MORLON, H., AND ETIENNE, R. S. The reconstructed tree in the lineage-based model of protracted speciation. *Journal of mathematical biology* 70, 1-2 (2015), 367–397.
- [43] LAMICHHANEY, S., BERGLUND, J., ALMÉN, M. S., MAQBOOL, K., GRABHERR, M., MARTINEZ-BARRIO, A., PROMEROVÁ, M., RUBIN, C.-J., WANG, C., ZAMANI, N., GRANT, B. R., GRANT, P. R., WEBSTER, M. T., AND ANDERSSON, L. Evolution of Darwin’s finches and their beaks revealed by genome sequencing. *Nature* 518, 7539 (2015), 371.
- [44] LIU, L., YU, L., AND EDWARDS, S. V. A maximum pseudo-likelihood approach for estimating species trees under the coalescent model. *BMC Evolutionary Biology* 10 (2010), 302.
- [45] LIU, S., LORENZEN, E., FUMAGALLI, M., LI, B., HARRIS, K., XIONG, Z., ZHOU, L., KORNELIUSSEN, T., SOMEL, M., BABBITT, C., WRAY, G., LI, J., HE, W., WANG, Z., FU, W., XIANG, X., MORGAN, C., DOHERTY, A., O’CONNELL, M., MCINERNEY, J., BORN, E., DALÉN, L., DIETZ, R., ORLANDO, L., SONNE, C., ZHANG, G., NIELSEN, R., WILLERSLEV, E., AND WANG, J. Population genomics reveal recent speciation and rapid evolutionary adaptation in polar bears. *Cell* 157, 4 (2014), 785–794.
- [46] LONG, C., AND KUBATKO, L. The effect of gene flow on coalescent-based species-tree inference. *Systematic Biology*.
- [47] MADDISON, W. P. Gene trees in species trees. *Systematic Biology* 46, 3 (1997), 523–536.
- [48] MALLET, J. Hybridization as an invasion of the genome. *Trends in ecology & evolution* 20, 5 (2005), 229–237.
- [49] MALLET, J. Hybridization, ecological races and the nature of species: empirical evidence for the ease of speciation. *Philosophical Transactions of the Royal Society B: Biological Sciences* 363, 1506 (2008), 2971–2986.
- [50] MALLET, J., BESANSKY, N., AND HAHN, M. W. How reticulated are species? *BioEssays* 38, 2 (2016), 140–149.
- [51] MARTIN, S. H., DASMAHAPATRA, K. K., NADEAU, N. J., SALAZAR, C., WALTERS, J. R., SIMPSON, F., BLAXTER, M., MANICA, A., MALLET, J., AND JIGGINS, C. D. Genome-wide evidence for speciation with gene flow in *Heliconius* butterflies. *Genome Research* 23, 11 (2013), 1817–1828.
- [52] MARTINSEN, G. D., WHITHAM, T. G., TUREK, R. J., AND KEIM, P. Hybrid populations selectively filter gene introgression between species. *Evolution* 55, 7 (2001), 1325–1335.

- [53] MAYR, E. *Systematics and the origin of species, from the viewpoint of a zoologist*. Harvard University Press, 1942.
- [54] MEIER, J. I., MARQUES, D. A., MWAIKO, S., WAGNER, C. E., EXCOFFIER, L., AND SEEHAUSEN, O. Ancient hybridization fuels rapid cichlid fish adaptive radiations. *Nature Communications* 8 (2017).
- [55] MILLER, W., SCHUSTER, S. C., WELCH, A. J., RATAN, A., BEDOYA-REINA, O. C., ZHAO, F., KIM, H. L., BURHANS, R. C., DRAUTZ, D. I., AND WITTEKINDT, N. E. Polar and brown bear genomes reveal ancient admixture and demographic footprints of past climate change. *Proceedings of the National Academy of Sciences* 109, 36 (2012), E2382–E2390.
- [56] MIRARAB, S., BAYZID, M. S., BOUSSAU, B., AND WARNOW, T. Statistical binning enables an accurate coalescent-based estimation of the avian tree. *Science* 346, 6215 (2014), 1250463.
- [57] MIRARAB, S., REAZ, R., BAYZID, M. S., ZIMMERMANN, T., SWENSON, M. S., AND WARNOW, T. ASTRAL: genome-scale coalescent-based species tree estimation. *Bioinformatics* 30, 17 (2014), i541–i548.
- [58] MORAN, P. A. P. Random processes in genetics. In *Mathematical Proceedings of the Cambridge Philosophical Society* (1958), vol. 54, Cambridge University Press, pp. 60–71.
- [59] MORJAN, C. L., AND RIESEBERG, L. H. How species evolve collectively: implications of gene flow and selection for the spread of advantageous alleles. *Molecular ecology* 13, 6 (2004), 1341–1356.
- [60] MORLON, H. Phylogenetic approaches for studying diversification. *Ecology Letters* 17, 4 (2014), 508–525.
- [61] NADEAU, N. J., MARTIN, S. H., KOZAK, K. M., SALAZAR, C., DASMAHAPATRA, K. K., DAVEY, J. W., BAXTER, S. W., BLAXTER, M. L., MALLET, J., AND JIGGINS, C. D. Genome-wide patterns of divergence and gene flow across a butterfly radiation. *Molecular Ecology* 22, 3 (2013), 814–826.
- [62] PAPADOPOULOU, A., BERGSTEN, J., FUJISAWA, T., MONAGHAN, M. T., BARRACLOUGH, T. G., AND VOGLER, A. P. Speciation and DNA barcodes: testing the effects of dispersal on the formation of discrete sequence clusters. *Philosophical Transactions of the Royal Society B: Biological Sciences* 363, 1506 (2008), 2987–2996.
- [63] PEASE, J. B., AND HAHN, M. W. Detection and polarization of introgression in a five-taxon phylogeny. *Systematic Biology* 64, 4 (2015), 651–662.

- 712 [64] PECCOUD, J., OLLIVIER, A., PLANTEGENEST, M., AND SIMON, J.-C. A continuum of genetic
713 divergence from sympatric host races to species in the pea aphid complex. *Proceedings of the*
714 *National Academy of Sciences* 106, 18 (2009), 7495–7500.
- 715 [65] PEREIRA, R. J., MONAHAN, W. B., AND WAKE, D. B. Predictors for reproductive isolation in a
716 ring species complex following genetic and ecological divergence. *BMC Evolutionary Biology* 11
717 (2011), 194.
- 718 [66] PINHO, C., AND HEY, J. Divergence with Gene Flow: Models and Data. *Annual Review of*
719 *Ecology, Evolution, and Systematics* 41, 1 (2010), 215–230.
- 720 [67] PYRON, R. A., AND BURBRINK, F. T. Phylogenetic estimates of speciation and extinction rates
721 for testing ecological and evolutionary hypotheses. *Trends in Ecology & Evolution* 28, 12 (2013),
722 729–736.
- 723 [68] ROSINDELL, J., CORNELL, S. J., HUBBELL, S. P., AND ETIENNE, R. S. Protracted speciation
724 revitalizes the neutral theory of biodiversity. *Ecology Letters* 13, 6 (2010), 716–727.
- 725 [69] RUNDLE, H. D., NAGEL, L., BOUGHMAN, J. W., AND SCHLUTER, D. Natural selection and
726 parallel speciation in sympatric sticklebacks. *Science* 287, 5451 (2000), 306–308.
- 727 [70] SAMADI, S., AND BARBEROUSSE, A. The tree, the network, and the species. *Biological Journal*
728 *of the Linnean Society* 89, 3 (2006), 509–521.
- 729 [71] SATO, A., O’HUGIN, C., FIGUEROA, F., GRANT, P. R., GRANT, B. R., TICHY, H., AND KLEIN, J.
730 Phylogeny of Darwin’s finches as revealed by mtDNA sequences. *Proceedings of the National*
731 *Academy of Sciences* 96, 9 (1999), 5101–5106.
- 732 [72] SCHLIEP, K. P. phangorn: phylogenetic analysis in R. *Bioinformatics* 27, 4 (2011), 592.
- 733 [73] SCHLUTER, D. Ecological causes of speciation. *Endless forms: species and speciation* (1998),
734 114–129.
- 735 [74] SEEHAUSEN, O. African cichlid fish: a model system in adaptive radiation research. *Proceedings*
736 *of the Royal Society of London B: Biological Sciences* 273, 1597 (2006), 1987–1998.
- 737 [75] SIPOS, B., MASSINGHAM, T., JORDAN, G. E., AND GOLDMAN, N. PhyloSim-Monte Carlo simu-
738 lation of sequence evolution in the R statistical computing environment. *BMC bioinformatics* 12,
739 1 (2011), 104.
- 740 [76] SLATKIN, M. Gene flow and the geographic structure of natural populations. *Sci-*
741 *ence(Washington)* 236, 4803 (1987), 787–792.

- 742 [77] SOLÍS-LEMUS, C., YANG, M., AND ANÉ, C. Inconsistency of species tree methods under gene
743 flow. *Systematic Biology* 65, 5 (2016), 843–851.
- 744 [78] SOUCY, S. M., HUANG, J., AND GOGARTEN, J. P. Horizontal gene transfer: building the web of
745 life. *Nature Reviews Genetics* 16, 8 (2015), 472–482.
- 746 [79] SOUSA-SANTOS, C., GANTE, H. F., ROBALO, J., CUNHA, P. P., MARTINS, A., ARRUDA, M.,
747 ALVES, M. J., AND ALMADA, V. Evolutionary history and population genetics of a cyprinid fish
748 (*Iberochondrostoma olisiponensis*) endangered by introgression from a more abundant relative.
749 *Conservation Genetics* 15, 3 (2014), 665–677.
- 750 [80] STADLER, T. Recovering speciation and extinction dynamics based on phylogenies. *Journal of*
751 *Evolutionary Biology* 26, 6, 1203–1219.
- 752 [81] STAMATAKIS, A. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large
753 phylogenies. *Bioinformatics* 30, 9 (2014), 1312–1313.
- 754 [82] TAMURA, K., BATTISTUZZI, F. U., BILLING-ROSS, P., MURILLO, O., FILIPSKI, A., AND KUMAR,
755 S. Estimating divergence times in large molecular phylogenies. *Proceedings of the National*
756 *Academy of Sciences* 109, 47 (2012), 19333–19338.
- 757 [83] WAHLBERG, N., WEINGARTNER, E., WARREN, A. D., AND NYLIN, S. Timing major con-
758 flict between mitochondrial and nuclear genes in species relationships of Polygoni butterflies
759 (Nymphalidae: Nymphalini). *BMC Evolutionary Biology* 9 (2009), 92.
- 760 [84] WILLIS, S. C., MACRANDER, J., FARIAS, I. P., AND ORTÍ, G. Simultaneous delimitation of
761 species and quantification of interspecific hybridization in Amazonian peacock cichlids (genus
762 *Cichla*) using multi-locus data. *BMC Evolutionary Biology* 12 (2012), 96.
- 763 [85] WU, C.-I. The genic view of the process of speciation. *Journal of Evolutionary Biology* 14, 6
764 (2001), 851–865.
- 765 [86] WU, Y. Coalescent-based species tree inference from gene tree topologies under incomplete
766 lineage sorting by maximum likelihood. *Evolution; international journal of organic evolution* 66, 3
767 (2012), 763–775.
- 768 [87] XU, B., AND YANG, Z. Challenges in species tree estimation under the multispecies coalescent
769 model. *Genetics* 204, 4 (2016), 1353–1368.
- 770 [88] YANG, Z. The BPP program for species tree estimation and species delimitation. *Current Zoology*
771 61, 5 (2015), 854–865.

- 772 [89] YU, Y., DONG, J., LIU, K. J., AND NAKHLEH, L. Maximum likelihood inference of reticulate
773 evolutionary histories. *Proceedings of the National Academy of Sciences* 111, 46 (2014), 16448–
774 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.

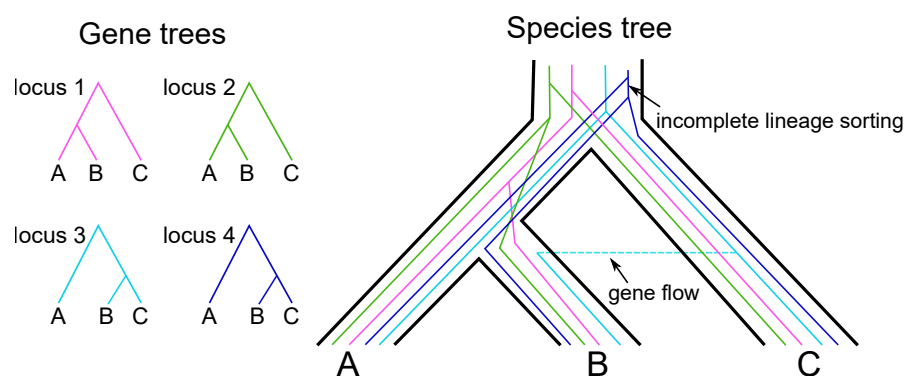


Figure 2: The genomic diversification (GD) model. Gene genealogies through species are depicted for two present-day genomes ($N = 2$ at $t = 0$) and four homologous genes ($n = 4$). Each gray ellipse represents a species (A-F). The model assumes that species are quasi-static in the timescale of a few generations, and each species lineage is located in a separate column. The genealogies of genes depend on four processes: introgression (*disconnection*), genetic differentiation (*intragenomic connection*), reproductive isolation (*intergenomic connection*), and replication (*coalescence*). In this example, the homologous genes 1 and 2 have coalesced.

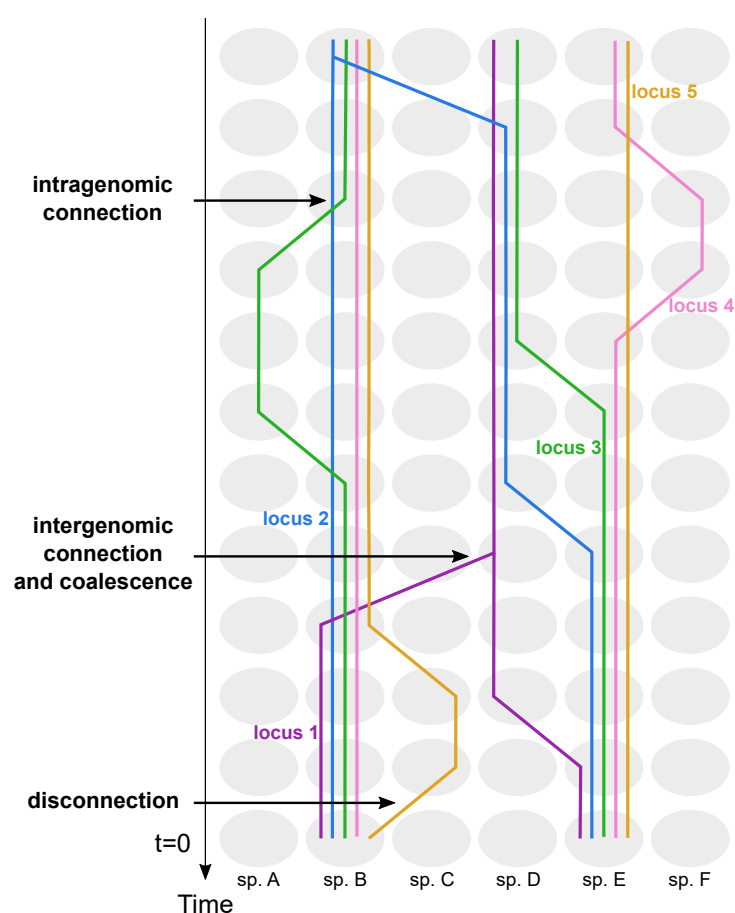


Figure 3: Evaluation of the GD model for a single sampled genome with n genes. Parameter settings: $a = 1$, $d \in [0.2, 2]$, every 0.2, and $n = 20, 50, 100$, and 200. The number of time units t was set to 10,000. We sampled the number of ancestral species every 500 time units starting at time $t = 5,000$, 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.

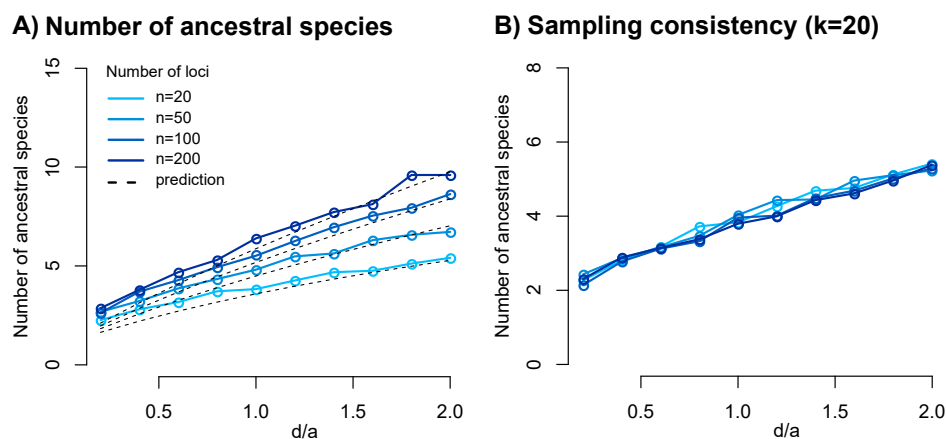


Figure 4: Billera-Holmes-Vogtmann (BHV) distances among sets of gene trees simulated under the genomic diversification (GD) model. For each set of parameters, 15 simulations were performed (with $t = 5,000$, enough to reach the coalescence of all homologous genes) and the median BHV distances were calculated. 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 *intergenomic connection* 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 $\frac{1}{a} = 0.3, 0.5, 0.9, 1.3, 1.7, 2.1, 2.5, 2.9, 3.3$.

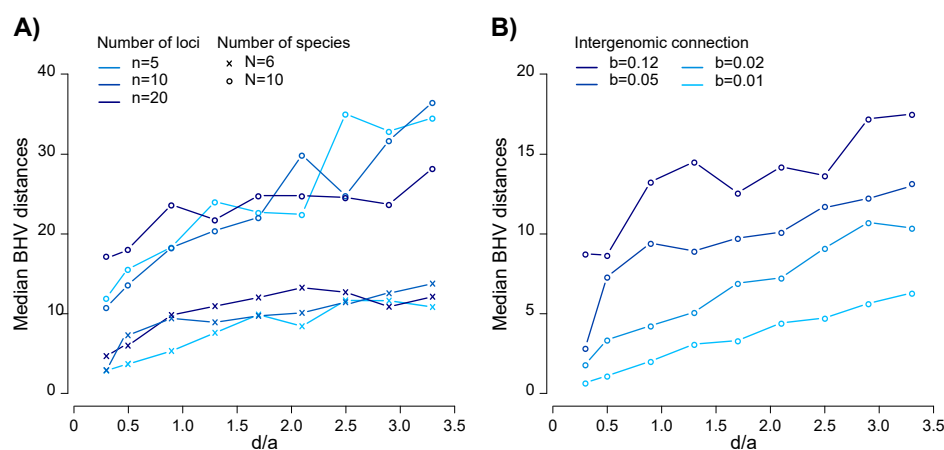


Figure 5: Bayesian phylogenetic reconstruction from simulated sequences under the GD model. We simulated 10 gene trees for 6 species under the GD model (with $\frac{b}{a} = 0.056$ and $\frac{d}{a} = 0.9$). The Bayesian phylogenetic analysis was performed with the program BEAST. The edges of the species tree (Bayesian analysis) are depicted by pipes in light gray. PP: posterior probabilities.

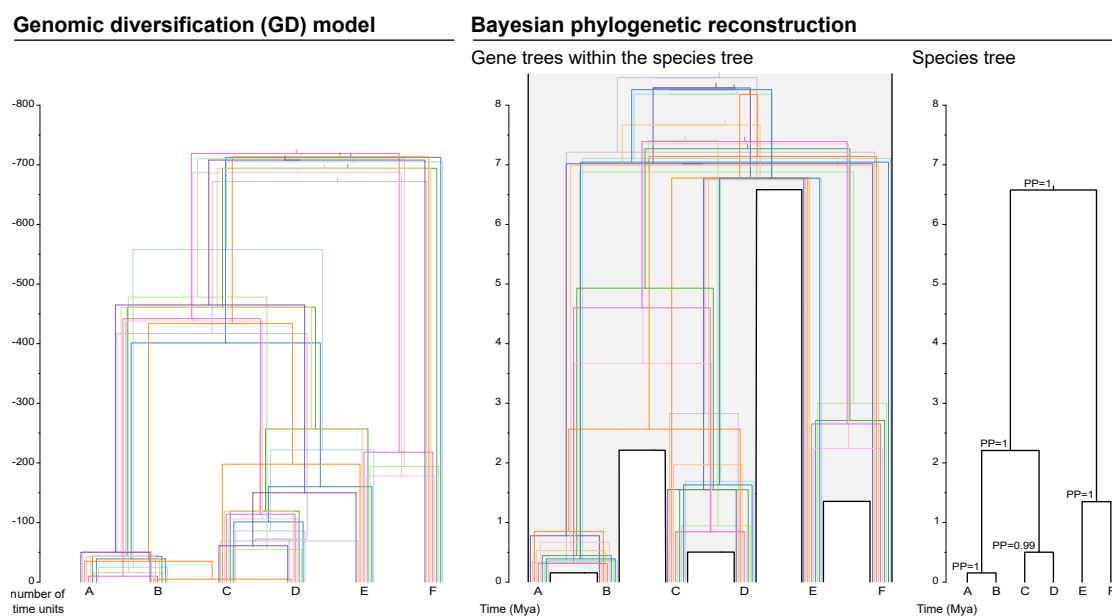
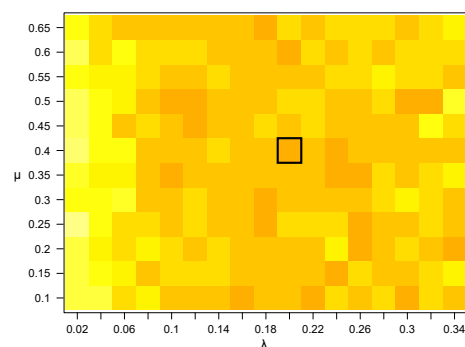


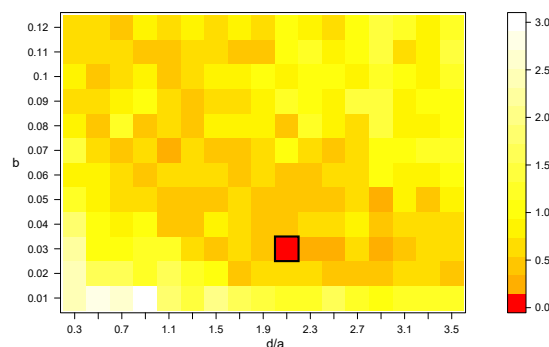
Figure 6: Minimization of the Kullback-Leibler (KL) divergence between empirical and simulated trees, *i.e.* between their distributions of BHV pairwise distances. Two parameters were optimized for each model. The *speciation* rate (λ) and the *extinction* rate (μ) for the multi-species coalescent (MSC) model (with coalescence rate set to 1). The *intergenomic connection* b and the ratio of the *disconnection* rate over the *intragenomic connection* rate ($\frac{d}{a}$) for the genomic diversification (GD) 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

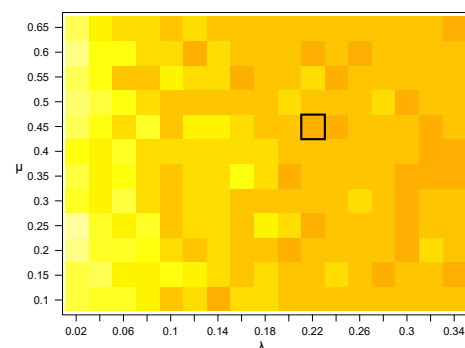


Genomic diversification (GD) model



Finch data-set

Multi-species coalescent (MSC) model



Genomic diversification (GD) model

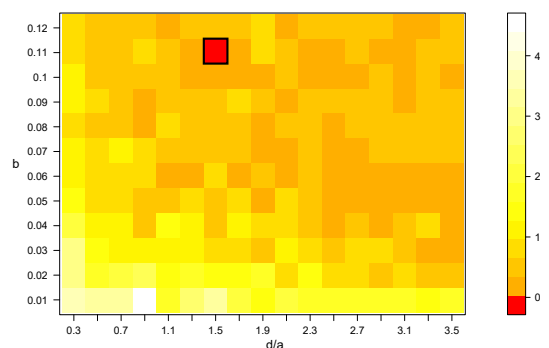
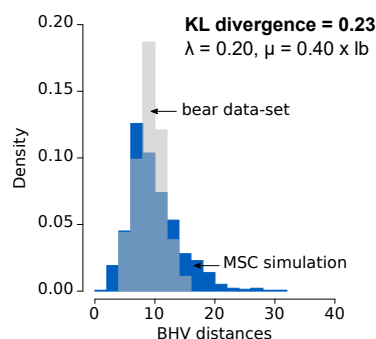


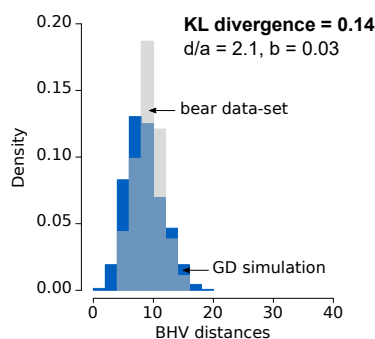
Figure 7: Best fit between empirical and simulated trees, *i.e.* between their distributions of BHV pairwise distances (selected cells of figure 6). For each set of variables, 15 simulations were performed and averaged. *a*: *intragenomic connection rate*, *b*: *intergenomic connection rate*, *d*: *disconnection rate* (set to 1), λ : *speciation rate*, μ : *extinction rate*, KL: Kullback-Leibler.

Bear data-set

Multi-species coalescent (MSC) model

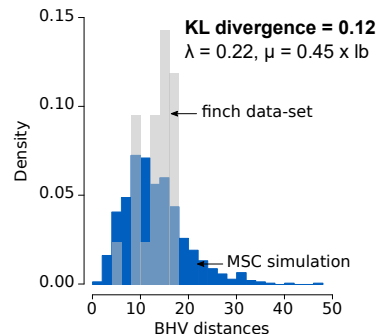


Genomic diversification (GD) model



Finch data-set

Multi-species coalescent (MSC) model



Genomic diversification (GD) model

