

POINTS OF VIEW

Potential survival of some, but not all, diversification methods

Brian C. O'Meara^{1*} and Jeremy M. Beaulieu²

¹Department of Ecology and Evolutionary Biology, University of Tennessee, Knoxville,
Tennessee, 37996-1610 USA

²Department of Biological Sciences, University of Arkansas, Fayetteville, Arkansas, 72701 USA

*Email for correspondence: bomeara@utk.edu

Abstract – Models have long been used for understanding changing diversification patterns over time. The rediscovery that models with very different rates through time can fit a phylogeny equally well has led to great concern about the use of these models. We share and add to these concerns: even with time heterogeneous models without these issues, the distribution of the data means that estimates will be very uncertain, something which is rarely communicated in empirical studies. Congruence issues such as those established for certain models of diversification also occur in models as basic as Brownian motion and coin flipping, and part of using models is learning when they can and cannot provide insights. The specific concern about lack of information about rates over time does not apply to models primarily seeking to understand rates across taxa (like in many uses of SSE models), but this does not prove immunity to incongruence in general in such models.

Keywords: diversification, congruence, likelihood, identifiability

For decades, molecular phylogenies have served as vital sources of historical information for deciphering the birth and the death dynamics of lineages. Hundreds, possibly thousands, of studies of molecular phylogenies have been dedicated to investigating diversification. In theory, if birth and death rates were indeed constant across the tree, estimating them separately is possible because each has distinguishable effects on the tree shape and branch length distributions (Nee et al. 1994). However, constant rates is likely far too simplistic of an assumption, and there are a number of extensions that expand this simple model for characterizing diversification as a function of time or diversity (e.g., Nee et al. 1992; Rabosky 2006, 2009; Bokma 2008; Rabosky and Lovette 2008; Morlon et al. 2011; Etienne et al. 2012), which are used to infer lines showing speciation and extinction rates scrolling into the past, like the pen of a seismometer tracking vibrations through time. A sudden sweep up of the extinction rate arm could mean a mass extinction. A slow, downward trajectory of the speciation rate arm as time approaches the present could mean available niches have become filled up, limiting the possibilities of adding new species. And, as with constant rate birth-death models, we have been working under the assumption that even tiny changes in speciation and/or extinction through time should leave distinct signatures on the tree shape and branching structure in a molecular phylogeny.

In a recent paper by Louca and Pennell (2020), the entire enterprise of estimating diversification rates, at least from molecular phylogenies alone, has been called into question (also see Pagel 2020). As it turns out, for any given phylogeny there are an infinite array of congruent models each having unique functions of speciation and/or extinction rates smoothly varying through time but with identical likelihood, and so indistinguishable from each other despite telling very different stories about the diversification history. This is based on the

property of both constant rate birth-death and time-varying models in which every lineage at any given time point experiences the same rates, and so sampling times for either a speciation or extinction event are drawn from the same distribution (also known as a coalescent point process or CPP; see Lambert and Stadler, 2013). Under such conditions, the likelihood of a tree under a given birth-death model can be inferred simply in terms of the lineage-through-time (LTT) curve, which is a retrospective counting of the number of lineages that led to a set of species observed today, and there are always multiple qualitatively different models that can produce the same curves with the same probability. For example, one model may infer the observed diversity of Cetaceans (i.e., whales, dolphins, and relatives) is a product of dramatic changes in the rate of speciation and extinction rates over time, whereas another, *equally likely* model, may infer modern whale diversity is the product of no extinction and ever so slight changes to the speciation rate. In other words, two diametrically opposed models, particularly with regards to the role of extinction, provide *equally* valid explanations for the mode and tempo of Cetacean diversification. In some cases, such as our example above, these models will have the same number of parameters, rendering them truly indistinguishable.

It should come as no surprise, then, that one popular interpretation of these findings is that any attempt to learn anything about diversification rates from molecular phylogenies is a completely futile enterprise. A different response, which we also have seen, is the continued and uncritical use of these suspect methods sanitized with a “but see Louca and Pennell (2020)” citation. It is also worth noting that the findings of Louca and Pennell (2020) are substantially similar, though much more detailed, to the ones presented by Kubo and Iwasa (1995) a quarter century ago. These authors also described an infinite array of birth and death models fitting the data equally well, which has been effectively ignored by most later workers.

The issues raised by Louca and Pennell (2020) and Kubo and Iwasa (1995) do represent substantial methodological problems for comparative biology. However, this does not signal the complete demise of studying diversification rates on molecular phylogenies, as some have claimed, as these problems do not extend to *all* models of diversification. Instead, they are limited to situations where the goal is to reconstruct and interpret diversification rates through time using what we refer to as, “time-varying, lineage homogeneous” models — again, models in which all lineages experience the same variable rates at any given point in time. These would be analogous to a non-heritable trait-dependent process (Lambert and Stadler, 2013), where changes in a trait occur exactly the same in all species independently (e.g., global CO₂, sea-level changes, global temperature patterns). Here, we briefly examine the case of time-constant, lineage heterogeneous (e.g., state-speciation and extinction, or SSE models; Maddison et al. 2007, Vasconcelos et al. 2022) models and show that they have access to information across clades not used by the former kind of models.

We also address some of the other procedures proposed, explicitly or implicitly, by Louca and Pennell (2020): continuing with pulled diversification rate reconstruction (or the “effective diversification rate” that includes the effect of sampling), focusing on a point estimate only, no longer penalizing for model complexity, and how information is distributed on trees.

Overall, we make four points:

1. Model congruence can occur in areas as different as coin flipping and Brownian motion: it does not mean these models must be given up, only that certain questions are infeasible.
2. Time-varying, lineage homogeneous models that use just the information from a lineage through time curve (even if the input is a full phylogeny) to estimate

changing speciation, extinction, diversification, turnover, or extinction fraction should be avoided due to congruence issues.

3. Pulled speciation and pulled diversification rate analyses (Louca and Pennell 2020) are identifiable, but they fail to incorporate the substantial uncertainty in reconstructions that come as a result of typically exponentially decreasing number of data points (lineages) as one approaches the root of a tree (this also plagues the methods in point 2)
4. Some SSE methods, and likely other methods that investigate heterogeneity across taxa, use information beyond that in a lineage through time curve and their utility remains intact in the face of Louca and Pennell (2020) and Kubo and Iwasa (1995).

Model congruence is common

It may come as a surprise that this issue of two models fitting data equally well is not new to comparative methods. Take, for instance, the inference of evolutionary trends, which, broadly defined, are identifiable patterns of trait evolution in a given direction through time. Using only extant species, can we detect horses getting bigger and with fewer digits, or increases in the mean seed size in flowering plants since the Cretaceous (e.g., Tifney 1984; Eriksson et al. 2000), or, more generally, uncover an evolutionary arms race between predator and prey (e.g., Dawkins and Krebs 1979; Abrams 1986)? It is trivial to extend a simple Brownian motion model to include a parameter that allows for the focal trait to evolve along a trend, and this is available in popular software like the R package *geiger* (Pennell et al. 2014). The likelihood for these models given the data is finite, and the simple no trend model is even nested within the trend model, so

comparisons between the two are straightforward. However, as Felsenstein (1988) and Hansen and Martins (1996) have pointed out, even though trait values move in a given direction under a Brownian motion with a trend model, this does not affect the expected covariances or means among species trait values. Consequently, the two models have identical likelihoods when fitted to extant species only, making them indistinguishable based on their probability alone. Careful biologists will not bother attempting to compare a trend versus no trend Brownian motion model, even if the question is compelling. Instead, we are limited to questions about some kinds of rates, not directions. For example, we know that the overall body size of Equidae has increased from the time of the group's origin, but using a phylogeny of modern horses and relatives we cannot infer this – we can ask how quickly lineages diverge in size from each other (the rate of variance accumulation over time) but this is not how quickly body size itself is changing unless the mean trend truly is zero.

However, with ancestral state reconstruction, it is easy to overlook the necessary assumptions. For example, the ancestor of a clade with body sizes ranging from 10-12 kg might have a reconstructed state near 11 kg under a model where the trend is forced to take an arbitrary value of zero but could have a reconstructed state of 50 kg under a model with a trend of an incremental trait decrease through time. By convention, we only use the former models, but the latter models could fit exactly as well. This makes ancestral state reconstruction dubious at best, but it does not mean that Brownian motion models are generally invalid for use on trees containing only modern taxa. We can still compare Brownian motion models with more complex models, such as Ornstein-Uhlenbeck models (e.g., Butler and King 2004; Beaulieu et al. 2012), Brownian models with more than one rate (e.g., O'Meara et al. 2006; Thomas et al. 2006), or models where the Brownian motion rate itself changes over time (e.g., Revell, 2021). In other

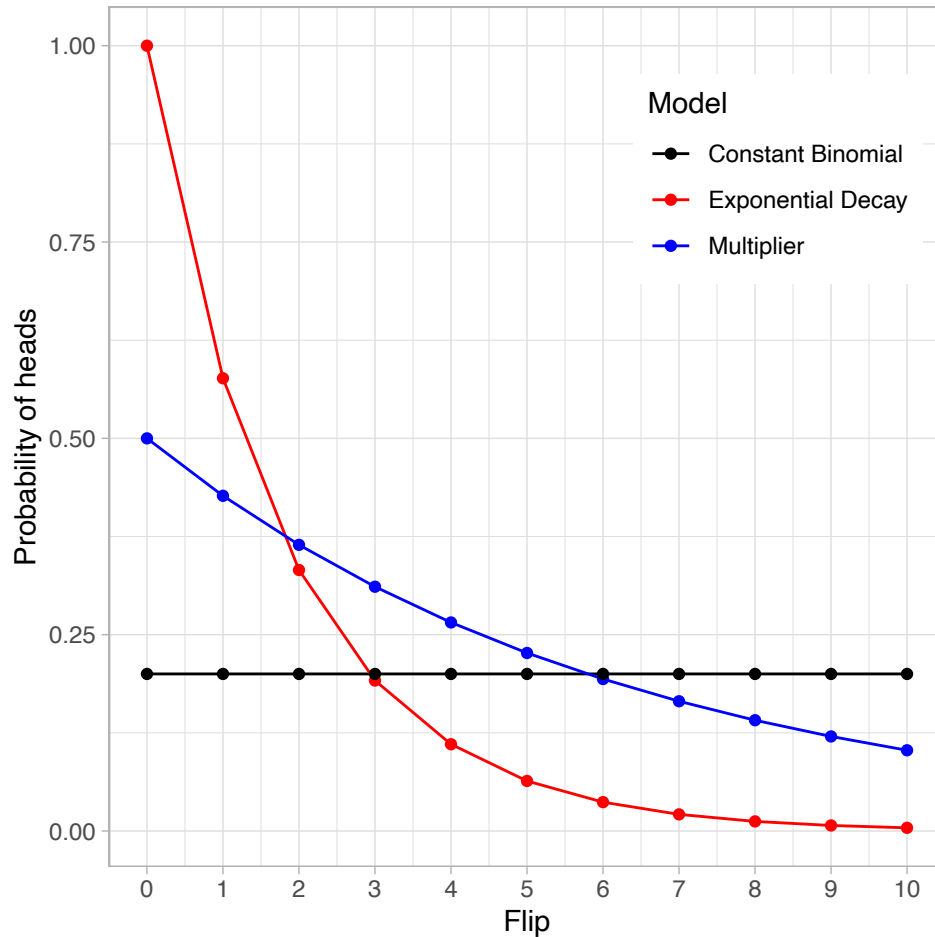


Figure 1. Probability of heads per flip on different models of coin flipping. Each of these models can fit the same dataset of two heads, eight tails with equal likelihood but make very different predictions about the next flip.

162

163 words, while Brownian motion with a trend model is unidentifiable with modern taxa only, we
 164 would not, for instance, say that any model that attempts to estimate rates of evolution on such
 165 trees is uninterpretable. Some models in this space give the same likelihoods and cannot be
 166 distinguished, but many others can, which calls for care and analysis, not panic. It does also point
 167 to issues that arise particularly when trying to estimate ancestral states.

168

We also point out that model congruence occurs in other statistically based disciplines.

169

Consider the classic coin-flipping example. Suppose we toss a coin 10 times, and 2 of those

SURVIVAL OF DIVERSIFICATION METHODS

tosses come up heads. The most straightforward fitted binomial model, which uses a single parameter, indicates that the probability of observing 2 heads in 10 flips is 0.3 for a biased coin with each flip having a 20% chance of landing on heads. Now suppose that every time we touch the coin, it gets slightly dented, or a bit of metal is worn away, and it becomes less and less likely to land on one side than the other. We can devise several models that have different slopes to alter the probability of heads after a set of coin flips (Figure 1). One model could allow every coin to start out with a 100% probability of getting heads and exponentially decrease at some rate down to zero probability of heads: this also has one parameter to estimate. Or one could start out with any coin being fair and its probability of heads being multiplied by a constant every flip, also with one parameter to measure. Interestingly, the probability of observing 2 heads in 10 flips of the coin in each of these models is the *same* as the simple binomial model, though the linear change models infer different initial probability of heads before any flips are made as well as what the probability of the next flip being heads. There are numerous other models, including using other distributions, that can fit a count of successes and failures: it will be hard or mathematically impossible to distinguish them. Despite all fitting the data in the present equally well, they may make different predictions about future flips and reconstructions of the probabilities of flips in the past. It does not mean coin flipping is impossible to model, nor that we cannot tell historical bias of a coin. It does mean that we are limited if we want to be able to infer how the probability of flips in the past has changed, and that we have tremendous uncertainty, beyond that estimated from any model, in predicting the results of future flips.

Avoid inference of congruent diversification models

While millions of students struggling with their statistics homework might cheer the destruction of the concept of estimating the probability of heads from a set of coin flips, it is important to emphasize that even though these models are functionally congruent and have the same number of parameters, many provide different predictions after a new set of coin flips are made (e.g., what is the likeliest outcome of the eleventh flip?). That is, even though they are indistinguishable from a probabilistic point of view, we can still likely distinguish them when new data becomes available. Of course, with comparative methods we cannot simply “flip” evolution more times to distinguish among a set of congruent models. The emphasis, then, as Morlon et al. (2022) recently pointed out, becomes what we are trying to learn about the world, given what we know about how it works. Incorporating additional information can also be important (Liow et al. 2022). For example, it is generally true that with coins, we have a good idea that the probability of heads does not change meaningfully over flips, so we may be willing to assume a standard binomial model and then question the fairness of a coin, perhaps as a way of extrapolating to other coins (i.e., if this Euro coin has a probability of heads of 0.502, is that true for other Euro coins?). In other words, the parameter can be of interest because the model is not really in question.

However, there is a risk in constraining to realistic models to ones limited by other sources of information (Liow et al. 2022). In other words, limiting ourselves to a perceived plausible region that may exclude the truth in the case of discoveries. A good example of this is the conflict over the age of the Earth between physicist Lord Kelvin and Darwin (Hattiangadi 1971). Kelvin had well-supported evidence that the Earth was less than 400 million years old, probably more like 200 million years old, and habitable for a subset of that. Darwin thought life

would take longer than that to have evolved its complexity, but his evidence was far less quantitative. Darwin turned out to be right. Kelvin did not know about the impact of radioactivity on the cooling of the Earth. It is a good thing that Darwin did not limit his ideas to those allowed by physics. On the other hand, if methods can only work well by being highly constrained by better sources of information, then they can be of little use. If dowsing for water only works in regions identified by using modern geological tools, then it is not clear what dowsing adds.

With many diversification models, the central question is about which model fits best, which is at odds with a general lack of knowledge about any system to clearly know which kind of model is appropriate ahead of time. In our view, we are not yet at the stage where we can confidently rule out a congruent model where extinction rates are driven by the position of a hypothetical dwarf star outside our solar system, which triggers periods of increased comet activity on Earth (e.g. Raup and Sepkowski 1984), over a more “sensible” model of, say, temperature clearly affecting speciation but not extinction rates. In such cases, asking questions about which of several indistinguishable models fit does not seem to us a good use of our time.

It is also important to emphasize that our argument here is not that the issues Kubo and Iwasa (1995) and Louca and Pennell (2020) point out are trivial. In fact, there are many papers, and even entire research programs, dedicated to the development of time-varying, lineage homogeneous models of diversification, and trying to draw conclusions based on which models fit best. But, as with coin flipping or Brownian motion, knowing what conclusions can be made given the models and data and limiting our work to those areas can be important. Moreover, if even coin flipping has congruent models, there is no guarantee that even models that currently seem to avoid the congruence issue, such as pulled diversification rates (the “effective diversification rate” that includes the effect of sampling) recommended by Louca and Pennell

(2020), or the many other methods that use trees to understand evolution, do not have other congruent models with different parameters, such as models that change rates by taxa rather than solely by time. Work on non-parametric identifiability (Stoult 2020; Louca and Pennell pers. comm.) may be fruitful for determining which models will end up being useful. As Box (1976) wrote, "Since all models are wrong the scientist must be alert to what is importantly wrong. It is inappropriate to be concerned about mice when there are tigers abroad." There could easily be a maladaptive radiation of papers showing how every model used in statistics, including those in biology, may have congruence issues or other failings – after all, if even coin flipping is suspect, what hope could there be for our far more complex domain? The key is to identify how the use or abandonment of such models gets us closer to understanding biology.

Avoid ancestral rate reconstruction

One area where abandonment seems prudent is the pursuit of most ancestral rate reconstructions. Ancestral state reconstruction of characters remains one of the most popular and widely used approaches in phylogenetic comparative methods, despite the occasional discussion to dampen enthusiasm in them (e.g., Cunningham et al. 1998; Omland 1999; Oakley and Cunningham, 2000). Ancestral state reconstruction is useful for formulating testable hypotheses, such as the synthesis and performance evaluation of putative ancestral proteins (e.g., Thornton et al. 2003; Pillai et al. 2020), biogeographic history and movements of clades through time (e.g., Ree and Smith, 2008; Landis et al. 2020), and the order and timing of character state changes (e.g., Schluter et al. 1997; Ackerly et al. 2006), though as mentioned above it can have issues in assuming no trend and other false assumptions. Reconstructing diversification rates through time has a similar appeal, in that they too can point to testable hypotheses about the intrinsic and

extrinsic factors that drive species diversity among groups. Armed with only a phylogeny of modern taxa, we can reconstruct how speciation rate, extinction rate, net diversification rate, or the new pulled diversification or pulled speciation rates, have changed through time. With the reconstruction of discrete or continuous characters, state information at the extant tips is generally less and less informative about states at nodes as one traverses deeper in the tree towards the root. For diversification rate models, the data are not arrayed along the tips of a tree, but rather, come from the distribution of branching events across the phylogeny. Ignoring uncertainty in branch lengths or topology, this makes a 10 Myr long edge equally informative regardless of whether it ended 3 million years ago or 30 million years ago.

In our view, many practitioners do not have a good intuitive sense of how information is distributed on a tree. Take, for example, Figure 2, which depicts a tree with one million taxa (from Louca and Pennell 2020). This tree gives 999,998 intervals between speciation events leading to extant tips, plus the interval after the last recovered speciation event, with which to estimate rates. The seemingly normal thing to do, which was done by Louca and Pennell (2020), is to split the tree into equal time bins (e.g., every 10 Myr) and estimate rates based on those bins. Even though the tree is far larger than any published study of diversification, they only estimate rates along 10 time intervals and for many of these bins there is only a trivial amount of data. For example, at the start of the 100 Myr to 90 Myr interval, there are just seven lineages, and by the end of that interval, there are only ten. The lineage through time plot, which is the data that goes into these methods, thus jumps just three times over that ten million years. This is clearly not a lot of data points for estimating speciation or extinction rates, or even a single pulled diversification rate. Each of the next several intervals have a single jump. That is, it goes from 10 to 11 lineages from 90 to 80 Myr, and from just 11 to 12 from 80 Myr to 90 Myr. It is no

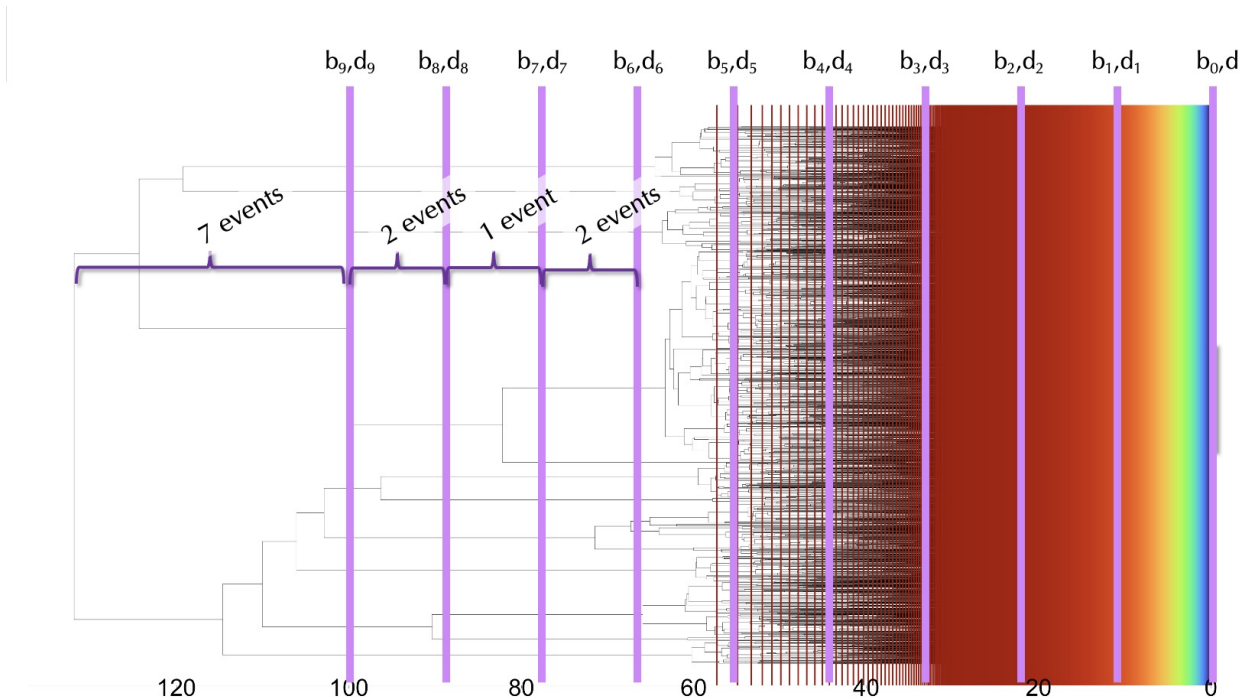


Figure 2. Million taxon tree from Louca and Pennell (2020). The purple lines separate the regimes used to estimate rates. The thin, rainbow-colored vertical lines separate regimes with 100 events within them representing equal-sized slices of data. Half the regimes are on each side of the green band, showing how much of the data are near the tips. The brackets show how many events occur in each regime.

285

286 wonder that these methods perform poorly; a single event on a 12 taxon tree does not contain

287 much information about rates, whether pulled or not. Put another way, these methods are

288 starving for data across large portions of the tree. Our intuition is that for big trees we have

289 information for much or all of their history, but in reality, nearly all the information is near the

290 tips. In fact, the midpoint of the data is the point at which half the number of lineages has

291 accumulated, which is not the halfway point along the time axis. Our failure of intuition comes

292 about, we feel, because we are not used to thinking of exponential branching processes (as

293 Maddison and FitzJohn 2015 noted, our field does not yet think in terms of the curvature of

294 biodiversity-time).

A natural corollary, then, is that reconstructions of the jiggles of rates backward in time (whether one does one rate per interval or allows a model to pick intervals) will contain increasing levels of uncertainty as one moves deeper in time. Nee et al. (1994) showed clearly that even rates from a constant birth-death model can carry substantial uncertainty. Yet most analyses doing the sort of work Louca and Pennell (2020) criticize, and even their examples, return a single point estimate for each parameter at a given time period. In a few cases, point estimates are summarized together across a set of trees, which is better, but still likely reflects substantially less uncertainty than what is truly present in any single estimate.

Besides unexamined uncertainty in point estimates, there is substantial uncertainty in which model fits best, even if one ignores the congruence issue. We took as an example of solid research work in the field, recent study by Condamine et al. (2020), which compared various models correlating various rates with angiosperm diversity (and other potential predictors) using just a phylogenetic tree; their best model showed an exponential dependence of conifer extinction rate with the number of angiosperms. This paper is far more careful than most of the genre, comparing fits of several models and using multiple realistic predictors. We tried adding more models to the set, including a model that fit random splines and ones where the predictors were time-scaled ratings of television shows. This is similar in spirit as Rabosky and Goldberg (2015) showing that whale names could provide a better fit of diversification rates than constant rate models. We wanted to see if biologically implausible predictors could result in good fits. In this paper's original model set, models nearly as good ($\Delta AICc < 2$; see their Table S5) include an effect on speciation or both speciation and extinction (only 41% of the model weight is on variable extinction only models; 39% is on variable speciation only, and 21% on both varying) – there is extremely little signal in the data on what exactly is varying, making it hard to draw

O'MEARA AND BEAULIEU

Conifer diversification with various predictors

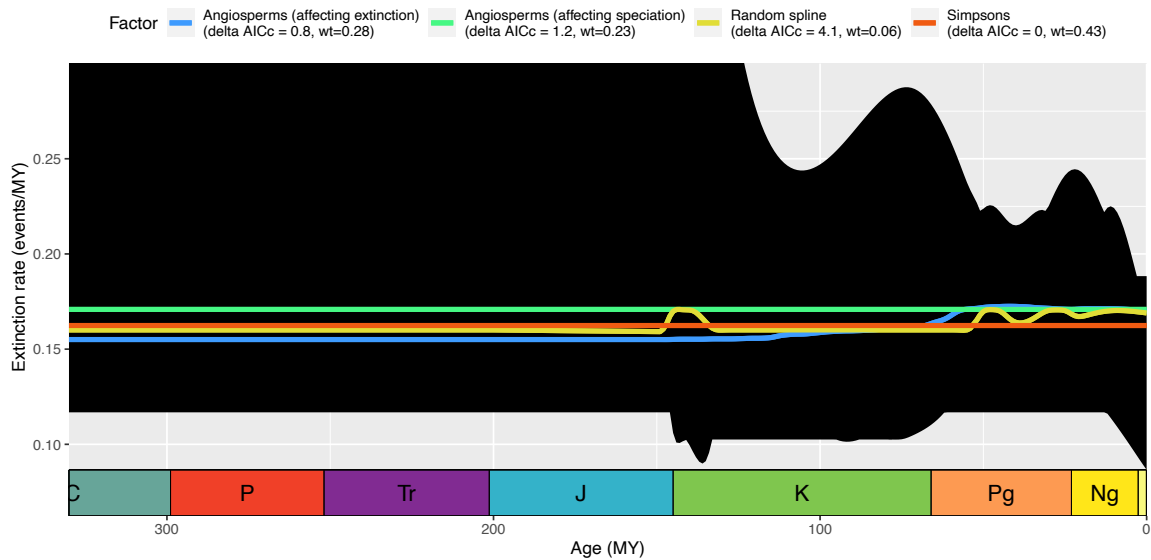
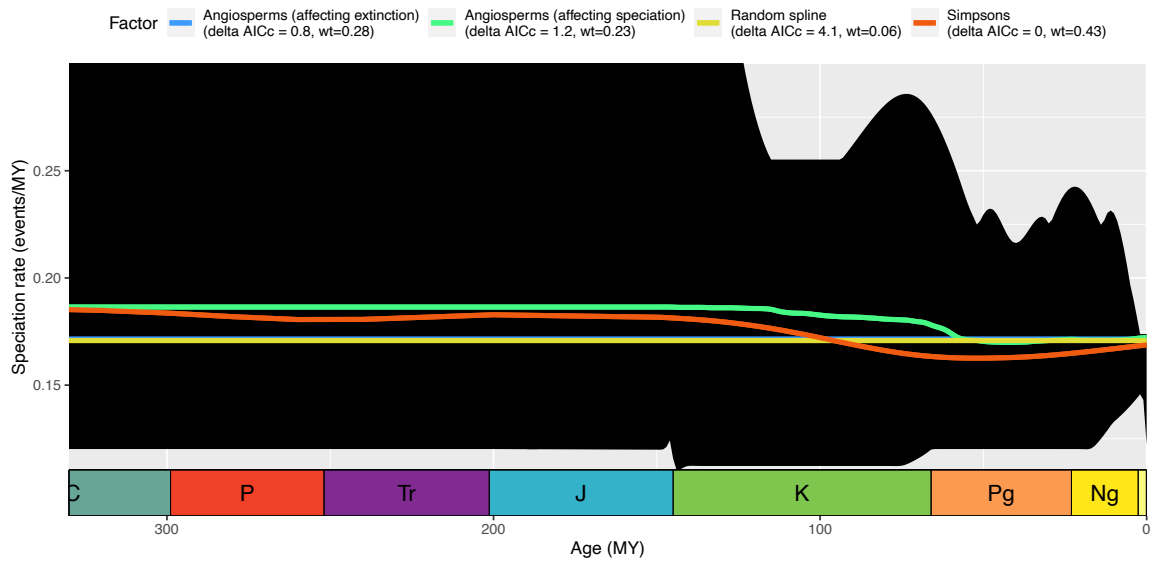
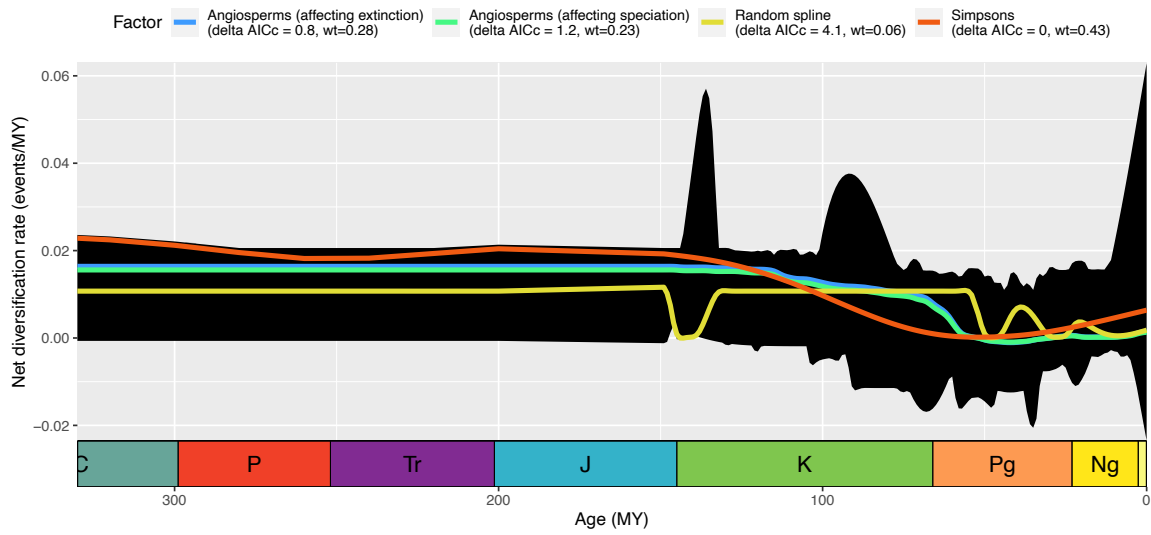


Figure 3. Comparison of net diversification, speciation, and extinction rate of conifers using as a predictor the best model from Condamine et al. (2020) in blue where only extinction rate varies with angiosperm diversity, a slightly worse model from that paper (green) where speciation rate varies with angiosperm diversity, a model (yellow) that fits the data best (at least in terms of likelihood — the number of free parameters of the spline is hard to compare), and using scaled IMDB ratings of the television program the Simpsons (red) as a predictor for speciation rate (which did a better job predicting conifer diversification than angiosperm diversity did). The black background shows the limits all the curves (generated with a variety of approaches) that were within 2.0 log likelihood units of the optimum, and so represents an estimate of the uncertainty of the rates. The true uncertainty is undoubtedly much higher; the need for new code to even attempt an estimate of the uncertainty for this figure points to a lack of attention to this issue in the field. Note that for legibility the speciation and extinction plots are truncated at a maximum rate of 0.3 events/MY, but the uncertainty goes far higher.

conclusions about mechanism from the best model. For example, by adding a random fitted (yellow) model in Figure 3, the diversification curves predict the conifer data even better but tell a very different story of constant speciation with decreases of extinction in the Cretaceous and Neogene rather than the recovered pattern of a gradual rise of extinction in the Cretaceous onward. Even using ratings of a television show (the Simpsons, the red line) scaled for the appropriate time period predicts conifer diversification better than the postulated angiosperm mechanism. This is not to say that we believe that conifers did have an extinction decrease in the Cretaceous and Neogene, nor that a television program at all relates to diversification. But even very careful work in this domain is left uncertain due to issues with these models.

Similarly, Morlon et al. (2011) looking at a paraphyletic set of 16 cetaceans found a constant speciation and variable extinction model fit best, but there were two other good models with a $\Delta AICc$ of less than 1 (including one where extinction does not vary) — this makes it hard to draw any firm conclusions from modern data alone. Careful biologists, as shown in the studies above, will limit themselves to only feasible mechanisms, but as we know from other diversification models (Rabosky and Goldberg 2015, Beaulieu and O'Meara, 2016), if presented with a very simple model and more complex alternatives only, methods using our messy,

complex empirical data will leap to use the more complex predictors. That is, if the only way to incorporate the very real heterogeneity of a process is to ascribe it to some varying predictor, methods will choose that. Whether it is 16 modern taxa or a million, it is unclear what we learn from such exercises. Our energies might be better directed elsewhere.

The state of SSE models and other approaches

Louca and Pennell (2020) speculate that state-speciation and extinction models (SSE) may have similar identifiability issues. This is not an unreasonable concern. Beaulieu and O'Meara (2016) demonstrated that if a trait has no effect on speciation and/or extinction rates, the likelihood of any SSE model becomes the product of the likelihoods of the Nee et al. (1994) tree likelihood and the character model likelihood (or the sum of the log-likelihoods in log space), so the models are clearly related. One could certainly alter the SSE model to include realistic factors like mass extinctions and secular changes in rates through time (we, and others, have experimented with these), and any one of these features will undoubtedly lead to a set of models with identical likelihoods. However, in other ways, strict SSE models can be immune to the particular concern of Louca and Pennell (2020), because they do not split the tree into time bins. Instead, they approximately treat a tree as a series of discrete chunks — that is, a chunk in one part of the tree is in state 0, and so is impacted by the instantaneous speciation rate, λ_0 , and extinction rate, μ_0 , while another chunk in another part of the tree is in state 1 and so is impacted by speciation rate, λ_1 , and extinction rate, μ_1 (in reality, they average over these paintings based on their probabilities). Within each of these chunks the speciation and extinction rates are invariant, and as Nee et al. (1994) showed, constrained in this way there is a single maximum likelihood estimate of each rate. If one limits the model space to where rates are dependent on

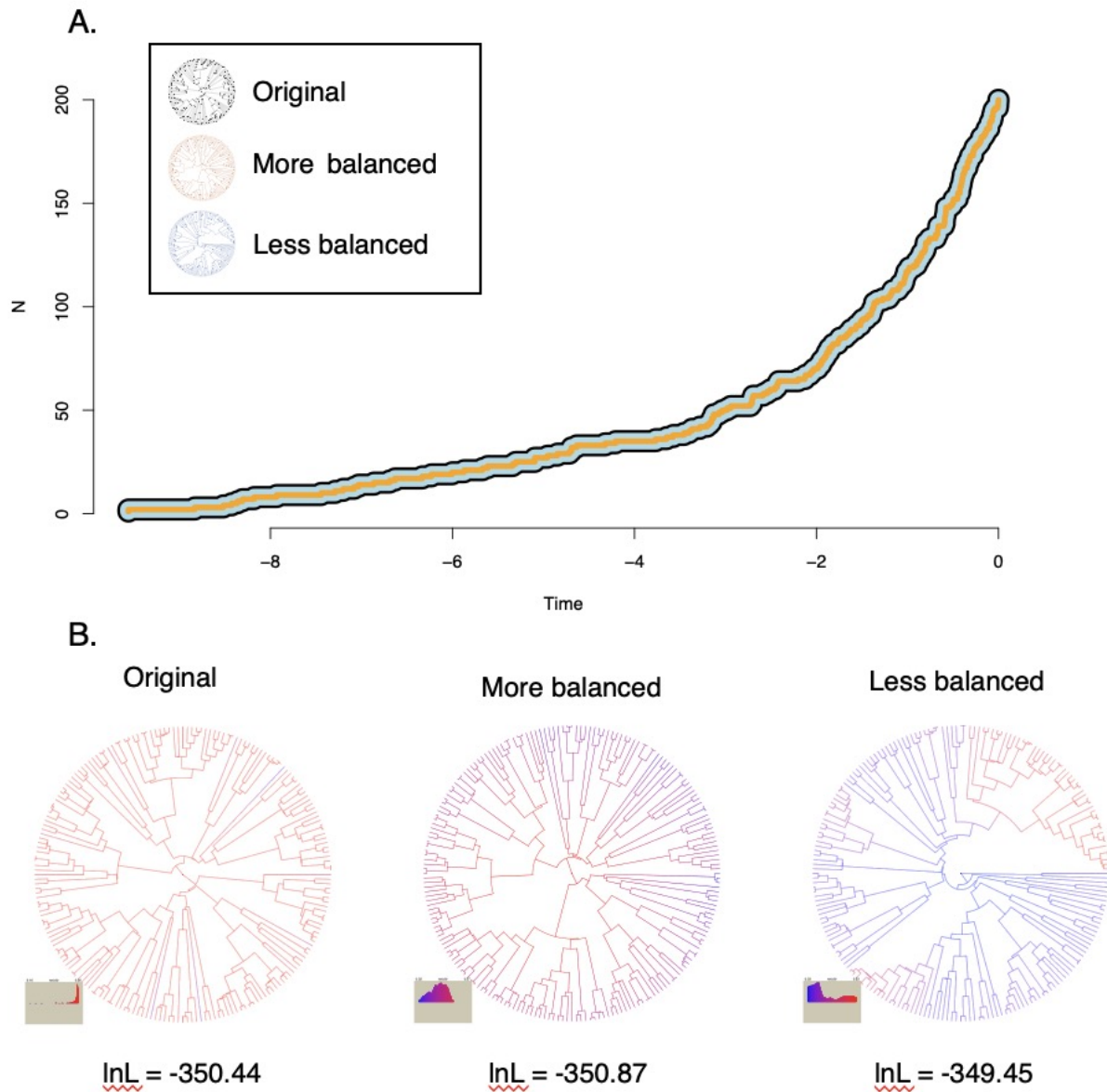


Figure 4: (A) Depicts the identical lineage through time (LTT) plots for three trees that differ in terms of tree balance. The procedure takes a simulated tree, then makes swaps across branches to either increase balance or decrease it but maintain the same lineage through time curve. (B) Depicts the log-likelihood score among the three trees under a two-rate MiSSE model. These trees produce identical log-likelihoods under taxon-homogeneous, time-heterogeneous models that use LTT data. However, this is not the case here because allowing rates to vary among clades, as our MiSSE models do, avoids the trap of having an infinite array of congruent models. Helmstetter et al. (2021) reach similarly positive conclusions about the possibility of learning about diversification from SSE models.

states (observed or hidden or some combination of both), then SSE models should be identifiable, though not immune to all the practical difficulties of estimating rates in the presence of extinction, finite data, errors in branch lengths and topology, and more.

We can at least empirically demonstrate that SSE models as constrained by current use are immune to the issues of model congruence based on information in the lineage through time plot: SSE models use more information than this. In Figure 4, there are three trees with identical lineage through time curves, but different arrangements of topology. Under a constant rate Yule or birth-death model the likelihoods of these three trees are identical, as one would expect given the findings of Louca and Pennell (2020). However, if we allow for multiple rates to be inferred across the tree by fitting a hidden states-only model (which we call MiSSE; see Vasconcelos et al. 2021) the three trees have different likelihoods. This is because the MiSSE model uses information not accessible to LTT methods, namely, the tree topology. Other methods that fit rate heterogeneity across taxa, such as BAMM (Rabosky 2014), MSBD (Barido-Sottani et al. 2018), and ClaDS (Maliot et al. 2019) may also not be bound by the issues that make different LTT models congruent. Even an approach as simple as sister group comparisons (e.g., Slowinski and Guyer 1993) can detect differences in the overall accumulation of species across pairs of clades in a way that depends on topology alone: identical lineage through time plots would have no effect on this. Taken together, this does not mean that clade-specific models of diversification could not have their own issues (even coin flipping models can have congruence, as shown above), just that the identifiability issue identified by Kubo and Iwasa (1995) and Louca and Pennell (2020) does not straightforwardly apply to them. Vasconcelos et al. (2022) showed that under a variety of complex diversification models, including scenarios involving multiple

regimes, diversity-dependent rates, clade-specific models with hidden rates can perform surprisingly well.

Still, there are scenarios that remain difficult for SSE models. For example, at the request of a reviewer, we also ran simulations of a model with no extinction and speciation rates exponentially increasing uniformly across the tree (simulated in the R package *castor*; Louca and Doebeli 2018). Thus, the generating model is outside the scope of typical SSE models in that it does not contain any state-specific or lineage-specific variation. When we analyzed these trees with MiSSE (Vasconcelos et al. 2022), despite the true underlying model having no extinction, MiSSE typically recovered an extinction rate about half the speciation rate. However, MiSSE did usually correctly find that the best model had no heterogeneity of rates. It also did a better job estimating speciation rates at the tips than net diversification rate (see Fig. S1). Nevertheless, given that both change continually through time it is a bit ambiguous what the “true” tip rate should be – for example, the rate at the instant the simulation ended, the average rate for the previous one million years, or some other weighted estimate. These tests demonstrate that when MiSSE does not find evidence for clade-specific rate variation, extracting qualitative meaning from speciation rates across the tips is still possible, but that estimates of extinction rates are likely to be suspect.

What are we really learning anyway?

Null hypothesis testing is intended to show whether an effect is significantly different from chance alone. At some point, though, comparing against chance becomes an uninteresting and dull exercise as the end point of a study. After several decades of studying diversification on molecular phylogenies and continually finding variation in rates across taxa and across time,

favoring a complex model over a “dull” null hypothesis of simple constant birth-death is no longer surprising. No reasonable scientist will argue that diversification processes have remained perfectly constant through time, with no changes in extinction rates, no factors changing speciation rates, and more. We know the data comes from a heterogeneous, complex process and so any complex model, even if somewhat reasonable, will fit better than a simple model. As we have noted elsewhere (see Beaulieu and O’Meara 2016; Caetano et al. 2018), rejecting the “null” does not imply that the slightly more complex alternative is the true model. Like a hot gas moved from a simple bottle to a more complex bottle with greater volume, our complex data will happily expand to take the shape of the biggest container offered to it. Model rejection, model weighting, posterior probability of models are all ways of saying, “my cloud of data is more comfortable in this larger bottle than in this smaller bottle. Since the extra bulge on the larger bottle is called factor X , this clearly shows that factor X is important.” However, a different bottle with the same volume but with a bulge for factor Y might fit as well. Good science will involve comparing different reasonable models to the data, not just comparing our slightly more complex model of interest with slightly simpler models. Much of our work on hidden rate models (e.g. Beaulieu et al. 2013; Beaulieu and O’Meara 2016; Caetano et al. 2018; Boyko and Beaulieu 2021) is motivated by this desire to give our preferred models an actual chance to lose against other models in the hope that we learn from this.

In our view, an important aspect of the work of Louca and Pennell (2020) was showing that even this limited, careful approach might not work for time-heterogeneous diversification rates: there are multiple diversification bottle shapes that fit the cloud of branching times from a tree equally well. Furthermore, approaches that seek to track the wiggles of diversification rates through time tell us very little, if anything, about the past. However, we would add that instead

of tracing the wiggles of a single pulled diversification rate, or even take the extreme step of stopping analyses of diversification using modern phylogenies altogether, we should use the valid methods we do have to answer biological questions, in the same way we can use Brownian motion even though different parameterizations can give identical likelihoods. Focus on analyses that lead to discoveries or confirmations of biological processes that are possible given available data.

On the whole, it is important to recognize that *our methods are better suited for using the past to learn about the present survivors, not using the present survivors to learn about the past.*

Phylogenies of extant taxa convey an enormous amount of information about species and their direct ancestors, but they also necessarily miss much of the history of a particular clade. As a consequence, there will never be a clever analysis of a phylogeny of extant archosaurs (crocodilians and birds) that will result in an inference of the dynamics of the rise and fall of sauropod dinosaurs, even though they are firmly nested in that clade and must have had a huge effect on the lineages that survived while all were interacting. Yet this is exactly what we are asking of our diversification analyses of modern taxa — that is, we think we are understanding something about diversification dynamics of archosaurs in the Cretaceous from a study of their weird, few surviving lineages. However, phylogenies of extant taxa can give us information about what led to present diversity, what traits are associated with modern diversity patterns, and, perhaps, even when certain modern lineages took off. We can understand something about diversification patterns of extant birds, for example, including what traits are associated with faster diversification or turnover rates. Moreover, even in cases where we have samples through time, model congruence can still be an issue (Louca et al. 2021).

Perhaps the best example of procedures that illustrate where we think the field needs to go are classic sister group comparisons (Mitter et al., 1988). These explicitly are about comparing modern clades and so are by their nature lineage-heterogeneous and limited to examining factors leading to modern diversity. They do not claim to allow inference about rate shifts in the past. There can be important corrections for even these methods (Käfer and Mousset, 2014) but they prevent scientists from spinning tales from limited information about the past. They should also be far more robust to the concerns raised by Maddison and FitzJohn (2015) than even hidden rate models. Of course, they are not without their own limitations: it can be hard to find enough comparisons; they only allow comparison of the direction of net diversification differences due to some pre-specified factor, while many of our hypotheses might relate to speciation rate, extinction rate, or, as we have advocated turnover rate (Beaulieu and O'Meara, 2016; Vasconcelos et al. 2021); they typically require only discrete characters (though see Harvey et al. 2020 and the bomeara/sisters package on GitHub); and they require ancestral state reconstruction to find sister pairs differing by a character state. There are also questions completely inaccessible to these methods; however, accepting these limitations at the outset may have prevented years of work that relied on methods that felt scientific but gave ultimately meaningless results given the issues now understood about time-heterogeneous diversification models.

Conclusions

The reconstruction of diversification rates through time, whether of pulled or classic rates, is appealing but flawed in the same way that inference of ancestral states is appealing but also flawed. Multiple indistinguishable models give very different estimates about the past, and

even for large trees, what matters is the branches and branching events at the times of interest, often when the mighty tree was a mere sapling. Moreover, this only looks at branches with modern descendants. What information it does provide is about what those lineages may have been doing, not what the clade as a whole may have been doing. Thus, approaches that seek to paint pictures about potential past diversification regimes at very incremental time periods are certainly suspect, with Louca and Pennell (2020) pointing to additional congruence issues that can affect diversification models.

We believe the best approach, given what we know now, is to avoid trying to estimate diversification rates through time from extant data. If practitioners continue to persist in this endeavor (and new metrics like pulled diversification rates have promise, with the caveat of difficulty in interpreting them), using normal statistical best practices are needed. Specifically, looking at multiple credible models, not relying on a single model for analysis when others are nearly as good, and paying attention to uncertainty in parameter estimates (and not just uncertainty from uncertainty in the tree, but the very substantial uncertainty present from estimating rates from mere handfuls of data). We can certainly learn about diversification processes from trees, but we need to recognize that what we can understand largely relates *only* to the surviving tips, with very little information on what happened along the way. Current SSE models and other models that infer rate heterogeneity across taxa, rather than across time, are more defensible in their use but still may have issues, albeit different from the particular issue raised by Kubo and Iwasa (1995) and Louca and Pennell (2020). Sister group analyses may grow in importance in future studies of factors that lead to different amounts of extant diversity.

Funding

This work was funded by the National Science Foundation grants DEB-1916558 and DEB-1916539.

Supplementary Materials

Code supporting this article is made freely available at <http://flippedcoin.info/>,
<https://github.com/bomeara/diversificationlives>, and
<https://github.com/bomeara/CondamineEtAlExample>.

Acknowledgements

We thank members of the Beaulieu and O'Meara labs for their comments and discussions of the ideas presented here. We would also like to thank Andrew Alverson, Jim Fordyce, and Ben Fitzpatrick for their insightful comments. Dan Rabosky and two anonymous reviewers provided useful feedback. This was funded by NSF grants DEB-1916539 to Brian O'Meara and DEB-1916558 to Jeremy Beaulieu.

References

- Abrams P.A. 1986. Is predator-prey coevolution an arms race? *Trends in Ecology and Evolution* 1:108-110.
- Ackerly D.D., Schilck D.W., Webb C.O. 2006. Niche evolution and adaptive radiation: testing the order of trait divergence. *Ecology* 87:S50-S61.
- Alexander H.K., Lambert A., Stadler T. 2016. Quantifying age-dependent extinction from species phylogenies. *Systematic Biology* 65:35-50.

- 518 Barido-Sottani J., Vaughan T.G., Stadler T. 2018. Detection of HIV transmission clusters from
519 phylogenetic trees using a multi-state birth-death model. *Journal of the Royal Society*
520 *Interface* 15:<https://doi.org/10.1098/rsif.2018.0512>
- 521 Beaulieu J.M., Jhwueng D.-C., Boettiger C., O'Meara B.C. 2012. Modeling stabilizing selection:
522 expanding the Ornstein-Uhlenbeck model of adaptive evolution. *Evolution* 66: 2369-
523 2383.
- 524 Beaulieu J.M., O'Meara B.C., Donoghue M.J. 2013. Identifying hidden rate changes in the
525 evolution of a binary morphological character: the evolution of plant habit in campanulid
526 angiosperms. *Systematic Biology* 62: 725:737.
- 527 Beaulieu J.M., O'Meara B.C. 2016. Detecting hidden diversification shifts in models of trait-
528 dependent speciation and extinction. *Systematic Biology* 65:583-601.
- 529 Boyko J.D., Beaulieu J.M. 2021. Generalized hidden Markov models for phylogenetic
530 comparative methods. *Methods in Ecology and Evolution* 12: 468:478.
- 531 Butler M.A., King A.A. 2004. Phylogenetic comparative analysis: a modeling approach for
532 adaptive evolution. *The American Naturalist* 164:683-695.
- 533 Bokma F. 2008. Bayesian estimation of speciation and extinction probabilities from (in)complete
534 phylogenies. *Evolution* 62:2441–2445.
- 535 Caetano D.S., O'Meara B.C., Beaulieu J.M. 2018. Hidden state models improve state-dependent
536 diversification approaches, including biogeographic models. *Evolution* 72:2308-2324.
- 537 Condamine F.L., Silvestro D., Koppelhus E.B., Antonelli A. 2020. The rise of angiosperms
538 pushed conifers to decline during global cooling. *Proceedings of the National Academy*
539 *of Sciences, USA* 117:28867-28875.

- 540 Condamine F.L., Rolland J., Morlon H. 2019. Assessing the causes of diversification
541 slowdowns: temperature-dependent diversity-dependent models receive equivalent
542 support. *Ecology Letters* 22:1900-1912.
- 543 Condamine F.L., Rolland J., Morlon H. 2013. Macroevoolutionary perspectives to environmental
544 change. *Ecology Letters* 16:72-85.
- 545 Cunningham C.W., Omland K.E., Oakley T.H. 1998. Reconstructing ancestral character states: a
546 critical reappraisal. *Trends in Ecology and Evolution* 13:361-366.
- 547 Dawkins R., Krebs J.R. 1979. Arms races between and within species. *Proceedings of the Royal*
548 *Society, B* 205:489-511.
- 549 Eriksson O., Friis E.M., Lofgren P. 2000. Seed size, fruit size, and dispersal systems in
550 angiosperms from the Early Cretaceous to the Late Tertiary. *The American Naturalist*
551 156:47-58.
- 552 Etienne R.S., Haegeman B., Stadler T., Aze T., Pearson P., Purvis A., Phillimore A. 2012.
553 Diversity-dependence brings molecular phylogenies closer to agreement with the fossil
554 record. *Proceedings of the Royal Society, B* 279:1300–1309.
- 555 Felsenstein J. 1988. Phylogenies and quantitative characters. *Annual Review of Ecology and*
556 *Systematics* 19:445-471.
- 557 Hansen T.F., Martins E.P. 1996. Translating between microevolutionary process and
558 macroevolutionary patterns: the correlation structure of interspecific data. *Evolution*
559 50:1404-1417.
- 560 Hattiangadi J.N. 1971. Alternatives and Incommensurables: The Case of Darwin and Kelvin.
561 *Philosophy of Science* 38:502-507.

- 562 Helmstetter A.J., Glemin S., Käfer J., Zenil-Ferguson R., Sauquet H., de Boer H., Dagallier M.J.,
563 Mazet N., Reboud E.L., Couvreur T.L.P., Condamine F.L. 2021. Pulled diversification
564 rates, lineages-through-time plots and modern macroevolutionary modeling. *Systematic*
565 *Biology* in press.
- 566 Käfer J., Mousset S. 2014. Standard sister clade comparison fails when testing derived character
567 states. *Systematic Biology* 63:601-609.
- 568 Kubo T., and Iwasa Y. 1995. Inferring the rates of branching and extinction from molecular
569 phylogenies. *Evolution* 49:694-704.
- 570 Lambert A., Stadler T. 2013. Birth-death models and coalescent point processes. *Theoretical*
571 *Population Biology* 90:113-128.
- 572 Landis M.J., Eaton D.A.R., Clement W.L., Park B., Spriggs E.L., Sweeney P.W., Edwards E.J.,
573 Donoghue M.J.. 2020. Joint phylogenetic estimation of geographic movements and
574 biome shifts during the global diversification of *Viburnum*. *Systematic Biology* 70:76–
575 94.
- 576 Liow L.H., Uyeda J., Hunt G.. 2022. Cross-disciplinary information for understanding
577 macroevolution. *Trends in Ecology and Evolution* *in press*
- 578 Louca S., Doebeli M. 2018. Efficient comparative phylogenetics on large trees. *Bioinformatics*
579 34:1053-1055.
- 580 Louca S., Pennell M.W. 2020. Extant timetrees are consistent with a myriad of diversification
581 histories. *Nature* 580:502-506.
- 582 Louca S., McLaughlin A., MacPherson A., Joy J.B., Pennell M.W.. 2021. Fundamental
583 identifiability limits in molecular epidemiology. *Molecular Biology and Evolution*
584 38:4010-4024.

- 585 Maddison W.P., FitzJohn R.G. 2015. The unsolved challenge to phylogenetic correlation tests
586 for categorical characters. *Systematic Biology* 64: 127-136.
- 587 Maddison W.P., Midford P.E., Otto S.P. 2007. Estimating a binary character's effect on
588 speciation and extinction. *Systematic Biology* 56:701–710.
- 589 Maliet O., Hartig F., Morlon H. 2019. A model with many small shifts for estimating species-
590 specific diversification rates. *Nature Ecology and Evolution* 3:1086-1092.
- 591 Marshall C. 2017. Five palaeobiological laws needed to understand the evolution of the living
592 biota. *Nature Ecology and Evolution*: 1: 1065
- 593 Mitter C.B., Farrell B., Wiegmann B. 1988. The phylogenetic study of adaptive zones: has
594 phytophagy promoted insect diversification? *American Naturalist* 132: 107-128.
- 595 Morlon H., Parsons T.L., Plotkin J.B. 2011. Reconciling molecular phylogenies with the fossil
596 record. *Proceedings of the National Academy of Sciences, USA* 108:16327–16332.
- 597 Morlon H., Robin S., Hartig F. 2023. Studying speciation and extinction dynamics from
598 phylogenies: addressing identifiability issues. *Trends in Ecology and Evolution* 37:497-
599 506.
- 600 Nee S., Holmes E. C., May R.M., Harvey P.H. 1994. Extinction rates can be estimated from
601 molecular phylogenies. *Philosophical Transactions of the Royal Society B* 344:77–82.
- 602 Nee S., Mooers A.Ø., Harvey P.H. 1992. The tempo and mode of evolution revealed from
603 molecular phylogenies. *Proceedings of the National Academy of Sciences, USA*
604 89:8322–8326.
- 605 Oakley T.H., Cunningham C.W.. 2000. Independent contrasts succeed where ancestor
606 reconstruction fails in a known bacteriophage phylogeny. *Evolution* 54:397-405.

SURVIVAL OF DIVERSIFICATION METHODS

- 607 O'Meara B.C., Ané C., Sanderson M.J., Wainwright P.C. 2006. Testing for different rates of
608 continuous trait evolution using likelihood. *Evolution* 60:922-933.
- 609 Omland K.E. 1999. The assumptions and challenges of ancestral state reconstructions.
610 *Systematic Biology* 48:604-611.
- 611 Pennell M.W., Eastman J.M., Slater G.J., Brown J.W., Uyeda J.C., FitzJohn R.G., Alfaro M.E.,
612 Harmon L.J. 2014. geiger v2.0: an expanded suite of methods for fitting
613 macroevolutionary models to phylogenetic trees. *Bioinformatics* 30:2216-2218.
- 614 Pillai A.S., Chandler S.A., Liu Y., Signore A.V., Cortez-Romero C.R., Benesch J.L.P.,
615 Laganowsky A., Storz J.F., Hochberg G.K.A., Thornton J.W. 2020. Origin of complexity
616 in haemoglobin evolution. *Nature* 581:480-485.
- 617 Rabosky D. L. 2006. Likelihood methods for detecting temporal shifts in diversification rates.
618 *Evolution* 60:1152–1164.
- 619 Rabosky D. L. 2009. Heritability of extinction rates links diversification patterns in molecular
620 phylogenies and fossils. *Systematic Biology* 58:629–640.
- 621 Rabosky D.L., Goldberg E.E. 2015. Model inadequacy and mistaken inferences of trait-
622 dependent speciation. *Systematic Biology* 64:340–355.
- 623 Rabosky D. L., Lovette I.J. 2008. Density-dependent diversification in North American wood
624 warblers. *Proceedings of the Royal Society, B* 275:2363–2371.
- 625 Raup D.M., Sepkowski, Jr, J.J. 1984. Periodicity of extinction in the geological past. *Proceedings*
626 *of the National Academy of Sciences, USA* 81:801-805.
- 627 Ree R.H., Smith S.A. 2008. Maximum likelihood inference of geographic range evolution by
628 dispersal, local evolution, and cladogenesis. *Systematic Biology* 57:4-14.

- 629 Revell L.J. 2021. A variable-rate quantitative trait evolution model using penalized likelihood.
630 bioRxiv doi: <https://doi.org/10.1101/2021.04.17.440282>.
- 631 Schluter D., Price T., Mooers A.Ø., Ludwig D. 1997. Likelihood of ancestor states in adaptive
632 radiation. *Evolution* 51:1699-1711.
- 633 Slowinski J.B., Guyer C. 1993. Testing whether certain traits have caused amplified
634 diversification: an improved method based on a model of random speciation and
635 extinction. *The American Naturalist* 142:1019-1024.
- 636 Stoult S. 2020. A Statistical Investigation of Species Distribution Models and Communication of
637 Statistics Across Disciplines. UC Berkeley. <https://escholarship.org/uc/item/2zq81799>
- 638 Thomas G.H., Freckleton R.P., Szekely T. 2006. Comparative analyses of the influence of
639 developmental mode on phenotypic diversification rates in shorebirds. *Proceedings of the*
640 *Royal Society, B* 273:1619-1624.
- 641 Thornton J.W., Need E., Crews D. 2003. Resurrecting the ancestral steroid receptor: ancient
642 origin of estrogen signaling. *Science* 301:1714-1717.
- 643 Tifney B.H. 1984. Seed size, dispersal syndromes, and the rise of angiosperms: evidence and
644 hypothesis. *Annals of the Missouri Botanical Garden* 71:551-576.
- 645 Vasconcelos T., O'Meara B.C, Beaulieu J.M. 2022. A flexible method for estimating tip
646 diversification rates across a range of speciation and extinction scenarios. *Evolution*
647 76:1420-1433.
- 648 Vasconcelos T. O'Meara B.C, Beaulieu J.M. 2022. Retiring “cradles” and “museums” of
649 biodiversity. *American Naturalist* 199:194-205.

650

651

Figure Legends

Figure 1: Probability of heads per flip on different models of coin flipping. Each of these models can fit the same dataset of two heads, eight tails with equal likelihood but make very different predictions about the next flip.

Figure 2: Million taxon tree from Louca and Pennell (2020). The purple lines separate the regimes used to estimate rates. The thin, rainbow-colored vertical lines separate regimes with 100 events within them representing equal-sized slices of data. Half the regimes are on each side of the green band, showing how much of the data are near the tips. The brackets show how many events occur in each regime.

Figure 3: Comparison of net diversification, speciation, and extinction rate of conifers using as a predictor the best model from Condamine et al. (2020) in blue where only extinction rate varies with angiosperm diversity, a slightly worse model from that paper (green) where speciation rate varies with angiosperm diversity, a model (yellow) that fits the data best (at least in terms of likelihood — the number of free parameters of the spline is hard to compare), and using scaled IMDB ratings of the television program the Simpsons (red) as a predictor for speciation rate (which did a better job predicting conifer diversification than angiosperm diversity did). The black background shows the limits all the curves (generated with a variety of approaches) that were within 2.0 log likelihood units of the optimum, and so represents an estimate of the uncertainty of the rates. The true uncertainty is undoubtedly much higher; the need for new code to even attempt an estimate of the uncertainty for this figure points to a lack of attention to this

issue in the field. Note that for legibility the speciation and extinction plots are truncated at a maximum rate of 0.3 events/MY, but the uncertainty goes far higher.

Figure 4: (A) Depicts the identical lineage through time (LTT) plots for three trees that differ in terms of tree balance. The procedure takes a simulated tree, then makes swaps across branches to either increase balance or decrease it but maintain the same lineage through time curve. (B) Depicts the log-likelihood score among the three trees under a two-rate MiSSE model. These trees produce identical log-likelihoods under taxon-homogeneous, time-heterogeneous models that use LTT data. However, this is not the case here because allowing rates to vary among clades, as our MiSSE models do, avoids the trap of having an infinite array of congruent models. Helmstetter et al. (2021) reach similarly positive conclusions about the possibility of learning about diversification from SSE models.