# Ranked Tree Shapes and the Future Loss of Phylogenetic Diversity

**Amaury Lambert** 

joint works with M. Steel, F. Gascuel and O. Maliet









June 26, 2018 Luminy

#### Outline

1. Loss of Phylogenetic Diversity

**2.** Introducing  $\beta$ 

3. Introducing Two Other Parameters : lpha and  $\eta$ 

4. Inference Results

► Q: "What fraction of the underlying evolutionary history survives when *k* of *n* species in a taxon are lost?" (Nee & May 1997)

- ► Q: "What fraction of the underlying evolutionary history survives when *k* of *n* species in a taxon are lost?" (Nee & May 1997)
- Phylogenies as metric trees, carry a footprint of evolutionary history

- ► Q: "What fraction of the underlying evolutionary history survives when *k* of *n* species in a taxon are lost?" (Nee & May 1997)
- Phylogenies as metric trees, carry a footprint of evolutionary history
- Phylogenetic Diversity PD = Total Length of Tree (Faith 1992)

- ► Q: "What fraction of the underlying evolutionary history survives when *k* of *n* species in a taxon are lost?" (Nee & May 1997)
- Phylogenies as metric trees, carry a footprint of evolutionary history
- Phylogenetic Diversity PD = Total Length of Tree (Faith 1992)
- Q becomes : "Can we predict how much Phylogenetic Diversity will remain in the face of present extinctions?"

### Field of Bullets Model (FoB)

- Take a phylogeny : tips = species
- ► Paint independently each tip in white w probability p, in black w probability 1 - p
- White dot = extant/sampled
- Black dot = extinct/not sampled



 Field of Bullets model : each species is removed independently, kept with probability p



(Mooers, Gascuel, Stadler, Li, Steel 2011)

- Field of Bullets model : each species is removed independently, kept with probability p
- Remaining PD S(p) = total length of tree spanned by extant/sampled species



(Mooers, Gascuel, Stadler, Li, Steel 2011)

- Field of Bullets model : each species is removed independently, kept with probability p
- Remaining PD S(p) = total length of tree spanned by extant/sampled species
- For a given tree, ES(p) is increasing and concave (Faller, Pardi, Steel 2008)

$$\mathbb{E}S(p) = \sum_{e} \ell(e) \left(1 - (1-p)^{n(e)}\right)$$

where :

 $\ell(e) =$ length of edge e, n(e) = # tips descending from e



(Mooers, Gascuel, Stadler, Li, Steel 2011)

#### Nee & May Science 1997

#### Extinction and the Loss of Evolutionary History

#### Sean Nee\* and Robert M. May

Extinction episodes, such as the anthropogenic one currently under way, result in a pruned tree of life. But what fraction of the underlying evolutionary history survives when k of n species in a taxon are lost? This is relevant both to how species loss has translated into a loss of evolutionary history and to assigning conservation priorities. Here it is shown that approximately 80 percent of the underlying tree of life can survive even when approximately 95 percent of species are lost, and that algorithms that maximize the amount of evolutionary history preserved are not much better than choosing the survivors at random. Given the political, economic, and social realities constraining conservation biology, these findings may be helpful.

### Nee & May Science 1997

"Approximately 80 percent of the underlying tree of life can survive even when approximately 95 percent of species are lost"

Rule of thumb :  $\mathbb{E}S_n(1) = \sum_{k=2}^n \frac{k}{\binom{k}{2}} \sim 2\log(n) \text{ so}$   $\frac{S_n(p)}{S_n(1)} \approx \frac{\log(pn)}{\log(n)} \approx 1$ 

- 1. Field of Bullets
- 2. Very Small External Edges
- 3. Balanced



#### The Kingman Coalescent

## Perfectly Balanced Tree (A) vs Caterpillar Tree (B)



## Loss of PD in Random Trees

Remaining PD is...

Low in imbalanced trees : more 'distinctive' sp

## Loss of PD in Random Trees

Remaining PD is...

- Low in imbalanced trees : more 'distinctive' sp
- ► High for the Kingman coalescent (Nee & May Science 1997)

## Loss of PD in Random Trees

Remaining PD is...

- Low in imbalanced trees : more 'distinctive' sp
- High for the Kingman coalescent (Nee & May Science 1997)
- Lower for the Yule tree (Mooers, Gascuel, Stadler, Li, Steel Syst Biol 2011) :

$$\frac{\mathbb{E}S(p)}{\mathbb{E}S(1)} = \text{ Ratio of expected remaining PD-to-Old PD } \approx -\frac{p\log p}{1-p}$$

#### Field of Bullets on a Birth-Death Tree

In a birth-death process stopped at time *T*, the **reduced tree** is a coalescent point process (CPP) : **node depths** are **i.i.d.** 



Lambert & Steel "Predicting the Loss of Phylogenetic Diversity under Non-Stationary Diversification Models" *JTB* 2013

CPP (e.g., reconstructed birth-death tree) : node depths H<sub>i</sub> are i.i.d.

Lambert & Steel "Predicting the Loss of Phylogenetic Diversity under Non-Stationary Diversification Models" *JTB* 2013

- CPP (e.g., reconstructed birth-death tree) : node depths H<sub>i</sub> are i.i.d.
- Conditional on *n* tips before FoB,

$$\lim_{n} \frac{S_{n}(p)}{S_{n}(1)} = p \frac{\mathbb{E}(B)}{\mathbb{E}(H)}$$

where

$$B:=\max_{i=1,\ldots,G}H_i,$$

and G is a geometric r.v. with success probability p.

Lambert & Steel "Predicting the Loss of Phylogenetic Diversity under Non-Stationary Diversification Models" *JTB* 2013

- CPP (e.g., reconstructed birth-death tree) : node depths H<sub>i</sub> are i.i.d.
- Conditional on *n* tips before FoB,

$$\lim_{n} \frac{S_{n}(p)}{S_{n}(1)} = p \frac{\mathbb{E}(B)}{\mathbb{E}(H)}$$

where

$$B:=\max_{i=1,\ldots,G}H_i,$$

and G is a geometric r.v. with success probability p.

Simple argument :

Lambert & Steel "Predicting the Loss of Phylogenetic Diversity under Non-Stationary Diversification Models" *JTB* 2013

- CPP (e.g., reconstructed birth-death tree) : node depths H<sub>i</sub> are i.i.d.
- Conditional on *n* tips **before** FoB,

$$\lim_{n} \frac{S_{n}(p)}{S_{n}(1)} = p \frac{\mathbb{E}(B)}{\mathbb{E}(H)}$$

where

$$B:=\max_{i=1,\ldots,G}H_i,$$

and G is a geometric r.v. with success probability p.

Simple argument :

► After FoB, the phylogenetic tree is a CPP with node depth *B* and *K<sub>n</sub>* tips

Lambert & Steel "Predicting the Loss of Phylogenetic Diversity under Non-Stationary Diversification Models" *JTB* 2013

- CPP (e.g., reconstructed birth-death tree) : node depths H<sub>i</sub> are i.i.d.
- Conditional on *n* tips **before** FoB,

$$\lim_{n} \frac{S_{n}(p)}{S_{n}(1)} = p \frac{\mathbb{E}(B)}{\mathbb{E}(H)}$$

where

$$B:=\max_{i=1,\ldots,G}H_i,$$

and G is a geometric r.v. with success probability p.

Simple argument :

- After FoB, the phylogenetic tree is a CPP with node depth B and K<sub>n</sub> tips
- By the SLLN,  $S_n(1) \sim n\mathbb{E}(H)$  and  $S_n(p) \sim K_n\mathbb{E}(B)$ ,

Lambert & Steel "Predicting the Loss of Phylogenetic Diversity under Non-Stationary Diversification Models" *JTB* 2013

- CPP (e.g., reconstructed birth-death tree) : node depths H<sub>i</sub> are i.i.d.
- Conditional on *n* tips **before** FoB,

$$\lim_{n} \frac{S_{n}(p)}{S_{n}(1)} = p \frac{\mathbb{E}(B)}{\mathbb{E}(H)}$$

where

$$B:=\max_{i=1,\ldots,G}H_i,$$

and G is a geometric r.v. with success probability p.

Simple argument :

- After FoB, the phylogenetic tree is a CPP with node depth B and K<sub>n</sub> tips
- By the SLLN,  $S_n(1) \sim n\mathbb{E}(H)$  and  $S_n(p) \sim K_n\mathbb{E}(B)$ ,
- Conclude with  $K_n/n \rightarrow p$ .

Lambert & Steel "Predicting the Loss of Phylogenetic Diversity under Non-Stationary Diversification Models" *JTB* 2013

For a Birth-Death tree with sp rate *b*, ext rate *d*, div rate r := b - d



Right : Slow progression towards the unit step function (from pure birth to critical) : d/b = 0 (the lowest curve) and then d/b = 0.5, 0.9, 0.99, 0.999.

- 1. What if trees are imbalanced?
  - $= \exists$  small clades = unprotected nodes

- 1. What if trees are imbalanced?
  - $= \exists$  small clades = unprotected nodes
- 2. What if older clades are more species-poor?
  - = deep nodes are in small clades

- 1. What if trees are imbalanced?
  - $= \exists$  small clades = unprotected nodes
- 2. What if older clades are more species-poor? = deep nodes are in small clades
- 3. What if older clades carry more extinct-prone species? = deep nodes are in first hit clades

- 1. What if trees are imbalanced?
  - $= \exists$  small clades = unprotected nodes
- 2. What if older clades are more species-poor? = deep nodes are in small clades
- 3. What if older clades carry more extinct-prone species? = deep nodes are in first hit clades

- 1. What if trees are imbalanced?
  - $= \exists$  small clades = unprotected nodes
- 2. What if older clades are more species-poor? = deep nodes are in small clades
- 3. What if older clades carry more extinct-prone species? = deep nodes are in first hit clades
- $\implies$  We need random tree models able to tune jointly :

- 1. What if trees are imbalanced?
  - $= \exists$  small clades = unprotected nodes
- 2. What if older clades are more species-poor? = deep nodes are in small clades
- 3. What if older clades carry more extinct-prone species? = deep nodes are in first hit clades
- $\implies$  We need random tree models able to tune jointly :
  - 1. Tree shape

 $\implies \beta =$  tree imbalance

- 1. What if trees are imbalanced?
  - $= \exists$  small clades = unprotected nodes
- 2. What if older clades are more species-poor? = deep nodes are in small clades
- 3. What if older clades carry more extinct-prone species? = deep nodes are in first hit clades
- $\implies$  We need random tree models able to tune jointly :
  - 1. Tree shape  $\implies \beta =$  tree imbalance
  - 2. Node depths' rankings

 $\implies \alpha =$ age-richness index

- 1. What if trees are imbalanced?
  - $= \exists$  small clades = unprotected nodes
- 2. What if older clades are more species-poor? = deep nodes are in small clades
- 3. What if older clades carry more extinct-prone species? = deep nodes are in first hit clades
- $\implies$  We need random tree models able to tune jointly :
  - 1. Tree shape  $\implies \beta =$  tree imbalance
  - 2. Node depths' rankings
  - 3. Abundances at tips

- $\implies \alpha =$ age-richness index
- $\implies \eta =$ abundance-richness index

#### Outline

1. Loss of Phylogenetic Diversity

#### 2. Introducing $\beta$

3. Introducing Two Other Parameters : lpha and  $\eta$ 

4. Inference Results

# Aldous' program

Aldous "Stochastic models and descriptive statistics for phylogenetic trees, from Yule to today" *Statist. Sci. 2001* 

1. What is a useful way to describe balance and imbalance in a general phylogenetic tree?

# Aldous' program

Aldous "Stochastic models and descriptive statistics for phylogenetic trees, from Yule to today" *Statist. Sci.* 2001

- 1. What is a useful way to describe balance and imbalance in a general phylogenetic tree?
- 2. Is there some particular region of the balance-imbalance spectrum containing most actual phylogenetic trees?

# Aldous' program

Aldous "Stochastic models and descriptive statistics for phylogenetic trees, from Yule to today" *Statist. Sci.* 2001

- 1. What is a useful way to describe balance and imbalance in a general phylogenetic tree?
- 2. Is there some particular region of the balance-imbalance spectrum containing most actual phylogenetic trees?
- 3. If so, is there some mathematically simple and biologically plausible stochastic model for phylogenetic trees whose realizations mimic actual trees?
• We are given distributions  $q_n$  on  $[n-1] = \{1, \ldots, n-1\}, \forall n \ge 2$ 

- We are given distributions  $q_n$  on  $[n-1] = \{1, \ldots, n-1\}, \forall n \ge 2$
- ▶ Recursively split each *k*-subset of [*n*] according to *q<sub>k</sub>* independently :

- We are given distributions  $q_n$  on  $[n-1] = \{1, \ldots, n-1\}, \forall n \ge 2$
- ▶ Recursively split each *k*-subset of [*n*] according to *q<sub>k</sub>* independently :

- We are given distributions  $q_n$  on  $[n-1] = \{1, \ldots, n-1\}, \forall n \ge 2$
- ▶ Recursively split each *k*-subset of [*n*] according to *q<sub>k</sub>* independently :



- We are given distributions  $q_n$  on  $[n-1] = \{1, \ldots, n-1\}, \forall n \ge 2$
- ▶ Recursively split each *k*-subset of [*n*] according to *q<sub>k</sub>* independently :



Induces a law on binary tree shapes with n labelled leaves.

- We are given distributions  $q_n$  on  $[n-1] = \{1, \ldots, n-1\}, \forall n \ge 2$
- ▶ Recursively split each *k*-subset of [*n*] according to *q<sub>k</sub>* independently :



- Induces a law on binary tree shapes with n labelled leaves.
- ► q<sub>n</sub> uniform yields the same tree shape as the reduced tree of any birth-death process (e.g., Yule tree)

A tree model is a family of probability distributions (P<sub>n</sub>) on (exchangeably labelled) tree shapes with n tips

• Call  $T_n$  a random tree with law  $P_n$ 

- Call T<sub>n</sub> a random tree with law P<sub>n</sub>
- ► Call  $T'_n$  the tree obtained by removing one tip from  $T_{n+1}$  (say the tip labelled n + 1)

- Call T<sub>n</sub> a random tree with law P<sub>n</sub>
- ► Call  $T'_n$  the tree obtained by removing one tip from  $T_{n+1}$  (say the tip labelled n + 1)
- ► The model is said **sampling consistent** if *T<sub>n</sub>* and *T'<sub>n</sub>* have the same distribution.

- Call T<sub>n</sub> a random tree with law P<sub>n</sub>
- ► Call  $T'_n$  the tree obtained by removing one tip from  $T_{n+1}$  (say the tip labelled n + 1)
- ► The model is said **sampling consistent** if *T<sub>n</sub>* and *T'<sub>n</sub>* have the same distribution.
- Example : Kingman coalescent.

## Aldous' Markov branching model (Mbm)

#### Theorem (Haas et al 2008, Lambert 2017)

A Mbm is sampling-consistent iff it there is a symmetric measure  $\nu$  on (0,1) s.t.

$$q_n(i) = a_n(\nu)^{-1} \binom{n}{i} \int_0^1 x^i (1-x)^{n-i} \nu(dx)$$

#### Construction

- Color dots are uniformly distributed in the interval
- Intervals are fragmented by r.v. R with law  $\sim \nu$  (red dashes)



Bertoin "Homogeneous fragmentation processes" PTRF 2001

A Markov process (Π<sub>t</sub>; t ≥ 0) with values in the partitions of N is a fragmentation process if it is non-increasing...

- A Markov process (Π<sub>t</sub>; t ≥ 0) with values in the partitions of N is a fragmentation process if it is non-increasing...
- ► ...and exchangeable = invariant by permutations (∃ block frequencies) + distinct blocks have independent futures.

- A Markov process (Π<sub>t</sub>; t ≥ 0) with values in the partitions of N is a fragmentation process if it is non-increasing...
- …and exchangeable = invariant by permutations (∃ block frequencies) + distinct blocks have independent futures.
- $\Pi$  is homogeneous if for any block *B* of  $\Pi(0)$  and for any bijection  $\varphi: B \to \mathbb{N}$ ,

$$(\varphi(\Pi_{|B}(t)); t \ge 0) \stackrel{\mathcal{L}}{=} (\Pi(t); t \ge 0).$$

Bertoin "Homogeneous fragmentation processes" PTRF 2001

- A Markov process (Π<sub>t</sub>; t ≥ 0) with values in the partitions of N is a fragmentation process if it is non-increasing...
- …and exchangeable = invariant by permutations (∃ block frequencies) + distinct blocks have independent futures.
- ►  $\Pi$  is homogeneous if for any block *B* of  $\Pi(0)$  and for any bijection  $\varphi: B \to \mathbb{N}$ ,

$$(\varphi(\Pi_{|B}(t)); t \ge 0) \stackrel{\mathcal{L}}{=} (\Pi(t); t \ge 0).$$

• Then  $\Pi_{[n]}$  is Markov  $\forall n$ 

- A Markov process (Π<sub>t</sub>; t ≥ 0) with values in the partitions of N is a fragmentation process if it is non-increasing...
- …and exchangeable = invariant by permutations (∃ block frequencies) + distinct blocks have independent futures.
- ►  $\Pi$  is homogeneous if for any block *B* of  $\Pi(0)$  and for any bijection  $\varphi: B \to \mathbb{N}$ ,

$$(\varphi(\Pi_{|B}(t)); t \ge 0) \stackrel{\mathcal{L}}{=} (\Pi(t); t \ge 0).$$

- Then  $\Pi_{[n]}$  is Markov  $\forall n$
- If Π is a binary, homogeneous fragmentation process, it is characterized by the fragmentation measure ν on (0, 1), where...

- A Markov process (Π<sub>t</sub>; t ≥ 0) with values in the partitions of N is a fragmentation process if it is non-increasing...
- …and exchangeable = invariant by permutations (∃ block frequencies) + distinct blocks have independent futures.
- $\Pi$  is homogeneous if for any block *B* of  $\Pi(0)$  and for any bijection  $\varphi: B \to \mathbb{N}$ ,

$$(\varphi(\Pi_{|B}(t)); t \ge 0) \stackrel{\mathcal{L}}{=} (\Pi(t); t \ge 0).$$

- Then  $\Pi_{[n]}$  is Markov  $\forall n$
- If Π is a binary, homogeneous fragmentation process, it is characterized by the fragmentation measure ν on (0, 1), where...
- ...each block fragments into two blocks with frequencies dx and 1 dx at rate  $\nu(dx)$ .

- A Markov process (Π<sub>t</sub>; t ≥ 0) with values in the partitions of N is a fragmentation process if it is non-increasing...
- …and exchangeable = invariant by permutations (∃ block frequencies) + distinct blocks have independent futures.
- $\Pi$  is homogeneous if for any block *B* of  $\Pi(0)$  and for any bijection  $\varphi: B \to \mathbb{N}$ ,

$$(\varphi(\Pi_{|B}(t)); t \ge 0) \stackrel{\mathcal{L}}{=} (\Pi(t); t \ge 0).$$

- Then  $\Pi_{[n]}$  is Markov  $\forall n$
- If Π is a binary, homogeneous fragmentation process, it is characterized by the fragmentation measure ν on (0,1), where...
- ...each block fragments into two blocks with frequencies dx and 1 dx at rate  $\nu(dx)$ .
- ► The previous theorem states the one-to-one correspondence : Sampling-consistent Mbm ⇐⇒ Π<sub>[[n]</sub> without fragmentation times

• The  $\beta$ -splitting model is for  $\beta \in (-2, \infty)$ :  $\nu(dx) = cx^{\beta}(1-x)^{\beta}dx$ 

- The  $\beta$ -splitting model is for  $\beta \in (-2, \infty)$ :  $\nu(dx) = cx^{\beta}(1-x)^{\beta}dx$
- Balance increases with  $\beta$

- The  $\beta$ -splitting model is for  $\beta \in (-2, \infty)$ :  $\nu(dx) = cx^{\beta}(1-x)^{\beta}dx$
- Balance increases with  $\beta$

• The  $\beta$ -splitting model is for  $\beta \in (-2, \infty)$ :  $\nu(dx) = cx^{\beta}(1-x)^{\beta}dx$ 

• Balance increases with  $\beta$ 



• The  $\beta$ -splitting model is for  $\beta \in (-2, \infty)$ :  $\nu(dx) = cx^{\beta}(1-x)^{\beta}dx$ 

• Balance increases with  $\beta$ 

Cat	PDA ERM	٨
-2	-1.5 -1 0	Beta
β	Description	Median split
-2	Completely unbalanced	1
-1.5	PDA model	1.5
-1	Unnamed	$\sqrt{m}$
0	Markov model	m/4
$\infty$	An almost completely balanced mod	el $m/2$

## Estimating $\beta$

 $S_{\min} VS S_{\min} + S_{\max}$  (Aldous 2001)

#### MLE of $\beta$ (Blum & François 2006)



My favorite evolutionary conundrum : Why  $\beta \approx -1$ ? (No answer here).

## Outline

1. Loss of Phylogenetic Diversity

**2.** Introducing  $\beta$ 

3. Introducing Two Other Parameters :  $\alpha$  and  $\eta$ 

4. Inference Results

Richness = Number of species in clade

Recall each subclade is associated w/ an interval of width X

- Recall each subclade is associated w/ an interval of width X
- Daughter subclades have widths  $X_{left} = RX$  and  $X_{right} = (1 R)X$

- Recall each subclade is associated w/ an interval of width X
- Daughter subclades have widths  $X_{left} = RX$  and  $X_{right} = (1 R)X$
- Iteratively **split nodes** from root to tips **w**/ **proba prop. to**  $\chi^{\alpha}$

- Recall each subclade is associated w/ an interval of width X
- Daughter subclades have widths  $X_{left} = RX$  and  $X_{right} = (1 R)X$
- Iteratively **split nodes** from root to tips **w**/ **proba prop. to**  $\chi^{\alpha}$
- Induces a sampling-consistent law on binary, ranked tree shapes with n labelled leaves.

- Recall each subclade is associated w/ an interval of width X
- Daughter subclades have widths  $X_{left} = RX$  and  $X_{right} = (1 R)X$
- Iteratively **split nodes** from root to tips **w**/ **proba prop. to**  $\chi^{\alpha}$
- Induces a sampling-consistent law on binary, ranked tree shapes with n labelled leaves.
- Yule tree/Kingman coalescent  $\Leftrightarrow \beta = 0$  and  $\alpha = 1$ .

- Recall each subclade is associated w/ an interval of width X
- Daughter subclades have widths  $X_{left} = RX$  and  $X_{right} = (1 R)X$
- Iteratively **split nodes** from root to tips **w**/ **proba prop. to**  $\chi^{\alpha}$
- Induces a sampling-consistent law on binary, ranked tree shapes with n labelled leaves.
- Yule tree/Kingman coalescent  $\Leftrightarrow \beta = 0$  and  $\alpha = 1$ .
- $\alpha > 0$  : deeper nodes in rich clades ('phylogenetic redundancy')

- Recall each subclade is associated w/ an interval of width X
- Daughter subclades have widths  $X_{left} = RX$  and  $X_{right} = (1 R)X$
- Iteratively **split nodes** from root to tips **w**/ **proba prop. to**  $\chi^{\alpha}$
- Induces a sampling-consistent law on binary, ranked tree shapes with n labelled leaves.
- Yule tree/Kingman coalescent  $\Leftrightarrow \beta = 0$  and  $\alpha = 1$ .
- $\alpha > 0$  : deeper nodes in rich clades ('phylogenetic redundancy')
- $\alpha < 0$  : deeper nodes in poor clades = PD at risk

- Recall each subclade is associated w/ an interval of width X
- Daughter subclades have widths  $X_{left} = RX$  and  $X_{right} = (1 R)X$
- Iteratively **split nodes** from root to tips **w**/ **proba prop. to**  $\chi^{\alpha}$
- Induces a sampling-consistent law on binary, ranked tree shapes with n labelled leaves.
- Yule tree/Kingman coalescent  $\Leftrightarrow \beta = 0$  and  $\alpha = 1$ .
- $\alpha > 0$  : deeper nodes in rich clades ('phylogenetic redundancy')
- $\alpha < 0$  : deeper nodes in poor clades = PD at risk
- See also : Sainudiin & Véber "A Beta-splitting model for evolutionary trees" Royal Society Open Science 2016



## Self-similar fragmentation processes

Bertoin "Self-similar fragmentation processes" Annales de l'IHP 2002

Recall the fragmentation process (Π(t); t ≥ 0) with values in the partitions of N.
Bertoin "Self-similar fragmentation processes" Annales de l'IHP 2002

- Recall the fragmentation process (Π(t); t ≥ 0) with values in the partitions of N.
- ▶  $\Pi^{\alpha}$  is self-similar with index  $\alpha$  if for any block *B* of  $\Pi^{\alpha}(0)$  with frequency *x* and for any bijection  $\varphi : B \to \mathbb{N}$ ,

$$(\varphi(\Pi^{\alpha}_{|B}(t));t\geq 0)\stackrel{\mathcal{L}}{=}(\Pi^{\alpha}(x^{\alpha}t);t\geq 0).$$

Bertoin "Self-similar fragmentation processes" Annales de l'IHP 2002

- Recall the fragmentation process (Π(t); t ≥ 0) with values in the partitions of N.
- ▶  $\Pi^{\alpha}$  is self-similar with index  $\alpha$  if for any block *B* of  $\Pi^{\alpha}(0)$  with frequency *x* and for any bijection  $\varphi : B \to \mathbb{N}$ ,

$$(\varphi(\Pi^{\alpha}_{|B}(t)); t \geq 0) \stackrel{\mathcal{L}}{=} (\Pi^{\alpha}(x^{\alpha}t); t \geq 0).$$

• Homogeneous fragmentation :  $\alpha = 0$ .

Bertoin "Self-similar fragmentation processes" Annales de l'IHP 2002

- Recall the fragmentation process (Π(t); t ≥ 0) with values in the partitions of N.
- ▶  $\Pi^{\alpha}$  is self-similar with index  $\alpha$  if for any block *B* of  $\Pi^{\alpha}(0)$  with frequency *x* and for any bijection  $\varphi : B \to \mathbb{N}$ ,

$$(\varphi(\Pi^{\alpha}_{|B}(t));t\geq 0)\stackrel{\mathcal{L}}{=}(\Pi^{\alpha}(x^{\alpha}t);t\geq 0).$$

- Homogeneous fragmentation :  $\alpha = 0$ .
- Now  $(\Pi^{\alpha}_{|[n]})$  is not Markov

Bertoin "Self-similar fragmentation processes" Annales de l'IHP 2002

- Recall the fragmentation process (Π(t); t ≥ 0) with values in the partitions of N.
- ▶  $\Pi^{\alpha}$  is self-similar with index  $\alpha$  if for any block *B* of  $\Pi^{\alpha}(0)$  with frequency *x* and for any bijection  $\varphi : B \to \mathbb{N}$ ,

$$(\varphi(\Pi^{\alpha}_{|B}(t));t\geq 0)\stackrel{\mathcal{L}}{=}(\Pi^{\alpha}(x^{\alpha}t);t\geq 0).$$

- Homogeneous fragmentation :  $\alpha = 0$ .
- Now  $(\Pi^{\alpha}_{|[n]})$  is not Markov
- But  $(\Pi_{|[n]}, R_n)$  is Markov, where  $R_n$  is the vector of *n*-tagged fragments = frequencies of blocks containing elements of [n]

SC Mbm w/ age-richness index  $\alpha \iff \prod_{l|n|}^{\alpha}$  w/ relative fragmentation times

Abundance = Relative abundance in number of individuals of clade

Each subclade assoc. w/ interval of width X, now also abundance A

- Each subclade assoc. w/ interval of width X, now also abundance A
- Daughter subclades have widths  $X_{left} = RX$  and  $X_{right} = (1 R)X$  and

$$A_{left} = \frac{|X_{left}|^{\eta}}{|X_{left}|^{\eta} + |X_{right}|^{\eta}} A = \frac{R^{\eta}}{R^{\eta} + (1 - R)^{\eta}} A$$
$$A_{right} = \frac{|X_{right}|^{\eta}}{|X_{left}|^{\eta} + |X_{right}|^{\eta}} A = \frac{(1 - R)^{\eta}}{R^{\eta} + (1 - R)^{\eta}} A$$

Abundance = Relative abundance in number of individuals of clade

- Each subclade assoc. w/ interval of width X, now also abundance A
- Daughter subclades have widths  $X_{left} = RX$  and  $X_{right} = (1 R)X$  and

$$A_{left} = \frac{|X_{left}|^{\eta}}{|X_{left}|^{\eta} + |X_{right}|^{\eta}} A = \frac{R^{\eta}}{R^{\eta} + (1 - R)^{\eta}} A$$
$$A_{right} = \frac{|X_{right}|^{\eta}}{|X_{left}|^{\eta} + |X_{right}|^{\eta}} A = \frac{(1 - R)^{\eta}}{R^{\eta} + (1 - R)^{\eta}} A$$

Induces a sampling-consistent law on binary tree shapes with (n labelled leaves and) a probability measure on its leaf set.

- Each subclade assoc. w/ interval of width X, now also abundance A
- Daughter subclades have widths  $X_{left} = RX$  and  $X_{right} = (1 R)X$  and

$$A_{left} = \frac{|X_{left}|^{\eta}}{|X_{left}|^{\eta} + |X_{right}|^{\eta}} A = \frac{R^{\eta}}{R^{\eta} + (1-R)^{\eta}} A$$
$$A_{right} = \frac{|X_{right}|^{\eta}}{|X_{left}|^{\eta} + |X_{right}|^{\eta}} A = \frac{(1-R)^{\eta}}{R^{\eta} + (1-R)^{\eta}} A$$

- Induces a sampling-consistent law on binary tree shapes with (n labelled leaves and) a probability measure on its leaf set.
- $\eta < 1$ : higher abundances in small clades

- Each subclade assoc. w/ interval of width X, now also abundance A
- Daughter subclades have widths  $X_{left} = RX$  and  $X_{right} = (1 R)X$  and

$$A_{left} = \frac{|X_{left}|^{\eta}}{|X_{left}|^{\eta} + |X_{right}|^{\eta}} A = \frac{R^{\eta}}{R^{\eta} + (1-R)^{\eta}} A$$
$$A_{right} = \frac{|X_{right}|^{\eta}}{|X_{left}|^{\eta} + |X_{right}|^{\eta}} A = \frac{(1-R)^{\eta}}{R^{\eta} + (1-R)^{\eta}} A$$

- Induces a sampling-consistent law on binary tree shapes with (n labelled leaves and) a probability measure on its leaf set.
- η < 1: higher abundances in small clades</p>
- $\eta = 1$ : total clade abundance prop. to richness

- Each subclade assoc. w/ interval of width X, now also abundance A
- Daughter subclades have widths  $X_{left} = RX$  and  $X_{right} = (1 R)X$  and

$$A_{left} = \frac{|X_{left}|^{\eta}}{|X_{left}|^{\eta} + |X_{right}|^{\eta}} A = \frac{R^{\eta}}{R^{\eta} + (1-R)^{\eta}} A$$
$$A_{right} = \frac{|X_{right}|^{\eta}}{|X_{left}|^{\eta} + |X_{right}|^{\eta}} A = \frac{(1-R)^{\eta}}{R^{\eta} + (1-R)^{\eta}} A$$

- Induces a sampling-consistent law on binary tree shapes with (n labelled leaves and) a probability measure on its leaf set.
- η < 1: higher abundances in small clades</p>
- $\eta = 1$ : total clade abundance prop. to richness
- $\eta > 1$ : lower abundances in small clades = PD at risk

- Each subclade assoc. w/ interval of width X, now also abundance A
- Daughter subclades have widths  $X_{left} = RX$  and  $X_{right} = (1 R)X$  and

$$A_{left} = \frac{|X_{left}|^{\eta}}{|X_{left}|^{\eta} + |X_{right}|^{\eta}} A = \frac{R^{\eta}}{R^{\eta} + (1-R)^{\eta}} A$$
$$A_{right} = \frac{|X_{right}|^{\eta}}{|X_{left}|^{\eta} + |X_{right}|^{\eta}} A = \frac{(1-R)^{\eta}}{R^{\eta} + (1-R)^{\eta}} A$$

- Induces a sampling-consistent law on binary tree shapes with (n labelled leaves and) a probability measure on its leaf set.
- η < 1: higher abundances in small clades</p>
- $\eta = 1$ : total clade abundance prop. to richness
- ▶  $\eta > 1$ : lower abundances in small clades = PD at risk
- Sampling (keeping extant) in prop to abundance ≈ Field of Bullets iff η = 1



•  $\eta = -3$ , rare species are all in large clades



▶  $\eta = -3$ , rare species are all in large clades

•  $\eta = 0$ , 'equal-sharing' measure (inv. prop. to number of splits)



▶  $\eta = -3$ , rare species are all in large clades

- $\eta = 0$ , 'equal-sharing' measure (inv. prop. to number of splits)
- ▶  $\eta = 1$ , species have equal abundances on average



▶  $\eta = -3$ , rare species are all in large clades

- ▶  $\eta = 0$ , 'equal-sharing' measure (inv. prop. to number of splits)
- ▶  $\eta = 1$ , species have equal abundances on average
- ▶  $\eta =$  3, the few species in poor clades are rare

### Warning!

The next slide represents 1 - S = PD loss...

...as a function of 1 - p = fraction of extinct species.

This function is expected to be increasing and convex in the FoB model.

# 'Danger zone' : $\alpha < \mathbf{0}, \eta > \mathbf{1}$



### Outline

1. Loss of Phylogenetic Diversity

**2.** Introducing  $\beta$ 

3. Introducing Two Other Parameters : lpha and  $\eta$ 

4. Inference Results

# Inference of $\alpha$



# Inference of $\eta$



### Data

#### Jetz et al "The global diversity of birds in space and time" Nature 2012



- ► Phylogenies of ≈ 100 clades in the class Aves (a.k.a. birds)
- Information used for inference :
  - Shapes of trees
  - Relative positions of nodes (but not the exact datation estimates)
  - Species range relative to sum of all species ranges.

### Inference from $\approx$ 100 bird clades



### Inference from $\approx$ 100 bird clades



α

# Conclusion

- New stochastic tree model with 3 parameters, tuning :
  - tree shape :  $\beta =$  tree balance
  - ranked node depths :  $\alpha = age$ -richness index
  - abundances at tips :  $\eta$  = abundance-richness index
- Danger zone :  $\beta < 0$ ,  $\alpha < 0$ ,  $\eta > 1$  = PD at risk
- Implemented in R-package apTreeshape (maintained by M.J. Blum)
- Large bird clades have  $\beta \approx -1$ ,  $\alpha < 0$ ,  $\eta \approx 1$
- Bird clades in danger zone!

Maliet, O., Gascuel, F., Lambert, A. (2018) "Ranked tree shapes, non-random extinctions and the loss of phylogenetic diversity" *Systematic Biology* (in press)

► To extend the fragmentation process  $\Pi$  to a random fragmentation triple  $(\Pi, W, M)$  (Jean-Jil Duchamps, work in progress), where :

- ► To extend the fragmentation process  $\Pi$  to a random fragmentation triple  $(\Pi, W, M)$  (Jean-Jil Duchamps, work in progress), where :
  - W(b) is the fragmentation rate of block b, instead of  $|b|^{\alpha}$

- ► To extend the fragmentation process  $\Pi$  to a random fragmentation triple  $(\Pi, W, M)$  (Jean-Jil Duchamps, work in progress), where :
  - W(b) is the fragmentation rate of block *b*, instead of  $|b|^{\alpha}$
  - M(b) is the mass measure of block *b*, instead of  $\sim |b|^{\eta}$

- ► To extend the fragmentation process  $\Pi$  to a random fragmentation triple  $(\Pi, W, M)$  (Jean-Jil Duchamps, work in progress), where :
  - W(b) is the fragmentation rate of block *b*, instead of  $|b|^{\alpha}$
  - M(b) is the mass measure of block *b*, instead of  $\sim |b|^{\eta}$
- ► To explain the triple conundrum :

- ► To extend the fragmentation process  $\Pi$  to a random fragmentation triple  $(\Pi, W, M)$  (Jean-Jil Duchamps, work in progress), where :
  - W(b) is the fragmentation rate of block *b*, instead of  $|b|^{\alpha}$
  - M(b) is the mass measure of block *b*, instead of  $\sim |b|^{\eta}$
- ► To explain the triple conundrum :
  - $\beta \approx -1$  : real trees are imbalanced

- ► To extend the fragmentation process  $\Pi$  to a random fragmentation triple  $(\Pi, W, M)$  (Jean-Jil Duchamps, work in progress), where :
  - W(b) is the fragmentation rate of block *b*, instead of  $|b|^{\alpha}$
  - M(b) is the mass measure of block *b*, instead of  $\sim |b|^{\eta}$
- ► To explain the triple conundrum :
  - $\beta \approx -1$  : real trees are imbalanced
  - $\alpha < 0$  : real trees have deeper nodes in smaller clades

- ► To extend the fragmentation process  $\Pi$  to a random fragmentation triple  $(\Pi, W, M)$  (Jean-Jil Duchamps, work in progress), where :
  - W(b) is the fragmentation rate of block *b*, instead of  $|b|^{\alpha}$
  - M(b) is the mass measure of block *b*, instead of  $\sim |b|^{\eta}$
- ► To explain the triple conundrum :
  - $\beta \approx -1$  : real trees are imbalanced
  - $\alpha < 0$  : real trees have deeper nodes in smaller clades
  - η ≥ 1: real trees have rarer species in smaller clades, but
    η ≈ 1: not much phylogenetic signal for species' abundances.

- ► To extend the fragmentation process  $\Pi$  to a random fragmentation triple  $(\Pi, W, M)$  (Jean-Jil Duchamps, work in progress), where :
  - W(b) is the fragmentation rate of block b, instead of  $|b|^{\alpha}$
  - M(b) is the mass measure of block *b*, instead of  $\sim |b|^{\eta}$
- ► To explain the triple conundrum :
  - $\beta \approx -1$  : real trees are imbalanced
  - $\alpha < 0$  : real trees have deeper nodes in smaller clades
  - $\eta \ge 1$ : real trees have rarer species in smaller clades, but  $\eta \approx 1$ : not much phylogenetic signal for species' abundances.
- ► To confirm the last two findings with more data.



Fanny Gascuel (Institut Curie) .....

Odile Maliet (ENS) .....





Mike Steel (U Canterbury) .....

SMILE : an interdisciplinary group in Paris









#### SMILE = Stochastic Models for the Inference of Life Evolution

# A positive answer to $\beta \approx -1$ ?

Hagen, Hartmann, Steel, Stadler "Age-Dependent Speciation Can Explain the Shape of Empirical Phylogenies" *Systematic Biology (2015)* 

► Birth-death process with age-dependent birth rate b = b(a) parameterized by  $b(a) = ca^{\phi-1}$ 





For  $\phi =$  0.6, the reconstructed tree has  $\beta \approx -$ 1.

Q: "Why  $\beta \approx -1$ ?"

# A positive answer to $\beta \approx -1$ ?

Hagen, Hartmann, Steel, Stadler "Age-Dependent Speciation Can Explain the Shape of Empirical Phylogenies" *Systematic Biology (2015)* 

► Birth-death process with age-dependent birth rate b = b(a) parameterized by  $\frac{3}{2}$ 

 $b(a)=ca^{\phi-1}$ 



 Estimates of \(\phi\) lie in (0, 1): speciation rate decreases with age

For  $\phi=$  0.6, the reconstructed tree has  $\beta\approx-$ 1.

Q: "Why  $\beta \approx -1$ ?"

- "Because  $\phi \approx$  0.6";-)
## Balance of incomplete trees



FIGURE 8. Effect of abundance-richness index  $\eta$  on the balance of phylogenetic trees after extinctions (Maximum Likelihood Estimate  $\hat{\beta}$  of  $\beta$ ). Initial tree balance  $\beta$  ranges from 10 (brown dots and lines, "bush trees") to -1.9 (green dots and lines, "caterpilar trees"). Extinction fraction pincreases from 0.01 to 0.98 (from left to right). Results are based on 100 simulation replicates: plain lines give median values and light areas give 9% confidence intervals. Other parameter values: number of species N = 100, approximation parameter  $\epsilon = 0.001$ ,  $\alpha = 0$ .