

# Shifted stochastic processes evolving on trees: application to models of adaptive evolution on phylogenies

Paul Bastide (1,2), Mahendra Mariadassou (2), Stéphane Robin (1)

(1) UMR 518 MIA, INRA/AgroParisTech, Paris, France.

(2) MaIAGE, INRA, Jouy-en-Josas, France.

## Goals and Setting

### Data

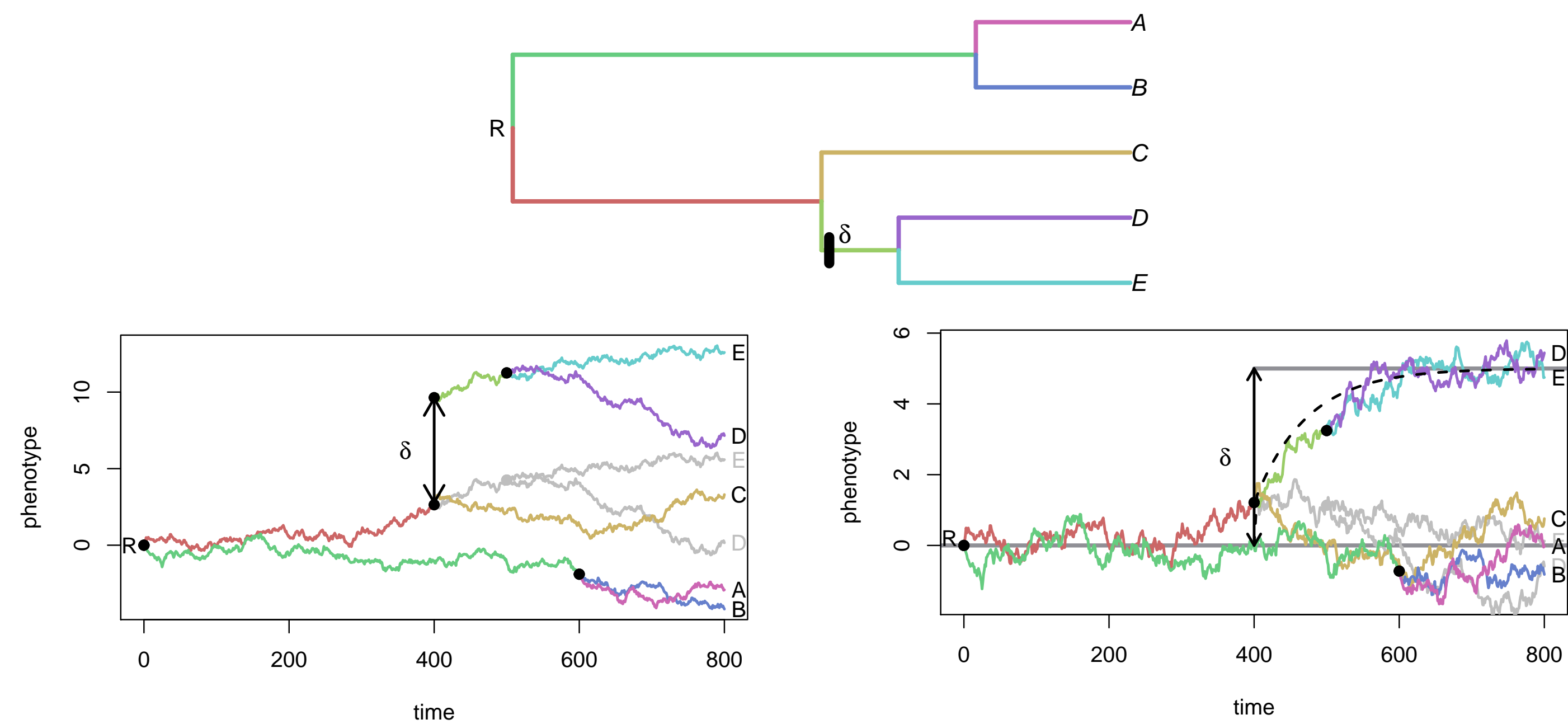
- Measure of one quantitative trait for a set of related extant species.
- A phylogenetic tree, time-calibrated and ultrametric.

### Goals

- Explain the observed trait distribution, while accounting for phylogenetic correlations.
- Detect environmental shifts that occurred in the past.

## Model

[2, 4]



$$dW(t) = \sigma dB(t)$$

BM Shifts in the **mean**.

$$dW(t) = \alpha[\beta(t) - W(t)]dt + \sigma dB(t)$$

OU Shifts in the **optimal value**.

## Incomplete Data Point of View

$$X_j | X_{pa(j)} \sim \mathcal{N} \left( q_j X_{pa(j)} + r_j + s_j \sum_k \mathbb{I}\{\tau_k = b_j\} \delta_k, \sigma_j^2 \right)$$

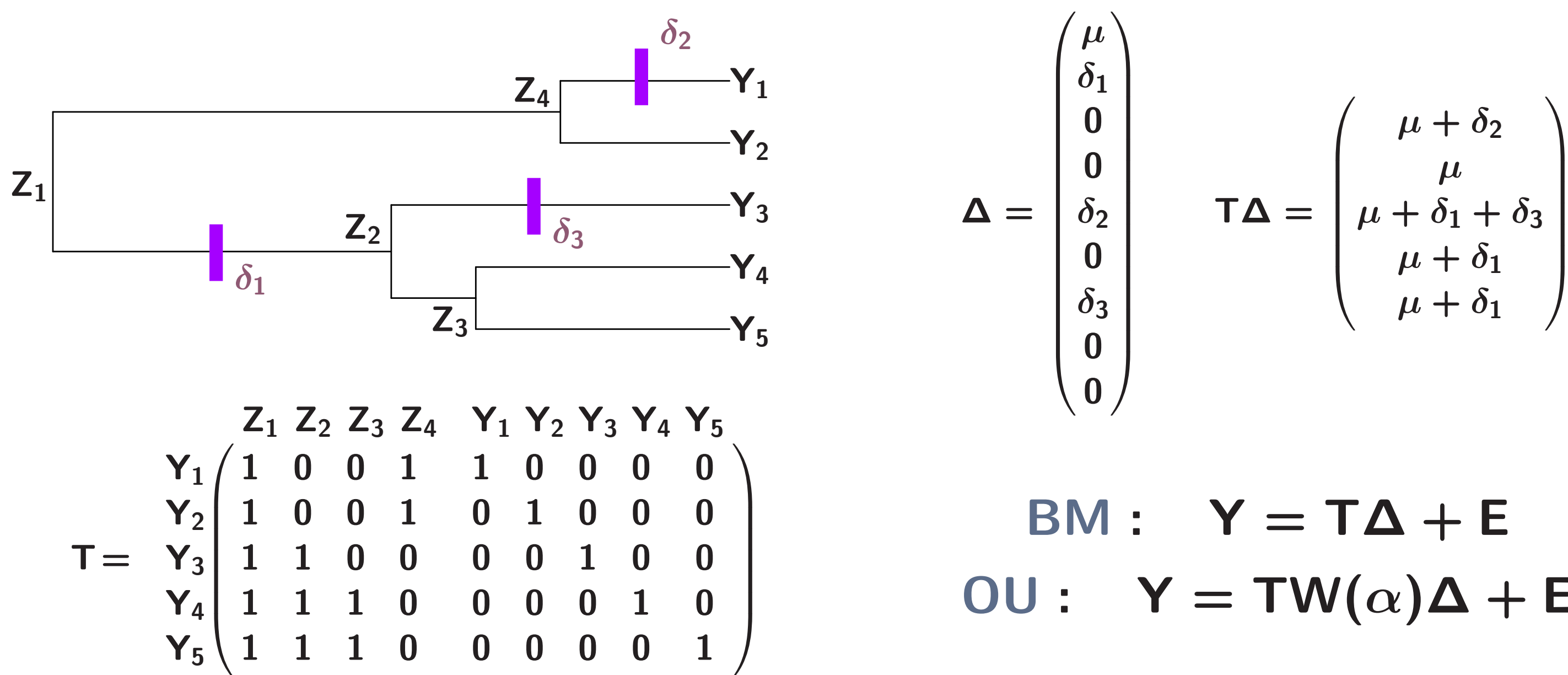
EM Algorithm Maximize  $\mathbb{E}_\theta[\log p_\theta(Z, Y) | Y]$ .

E step "Upward-Downward" Algorithm.

M step OU: increase objective function (GM).

Initialization LASSO regression.

## Linear Regression Point of View

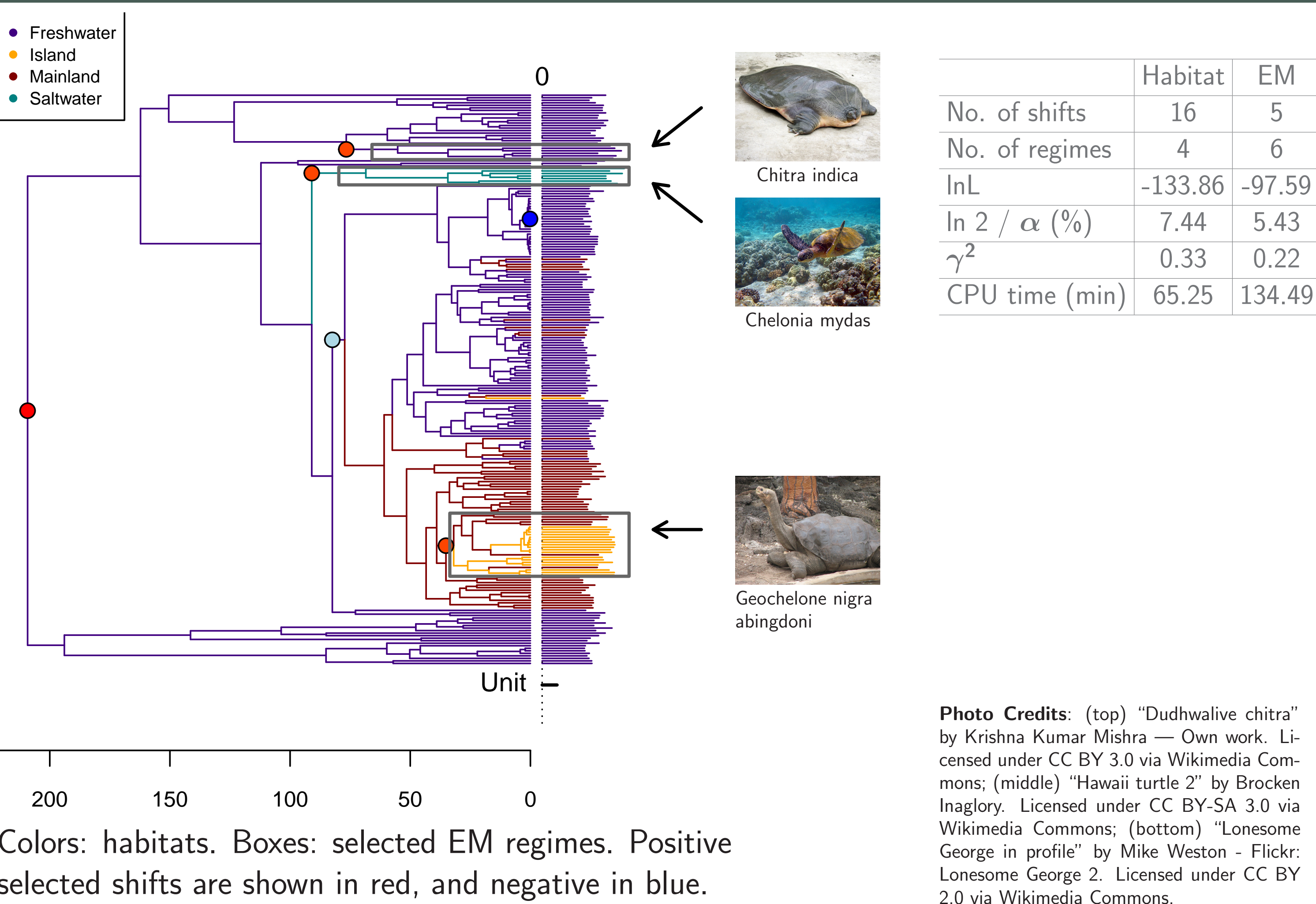


$$\text{BM : } Y = T\Delta + E$$

$$\text{OU : } Y = TW(\alpha)\Delta + E$$

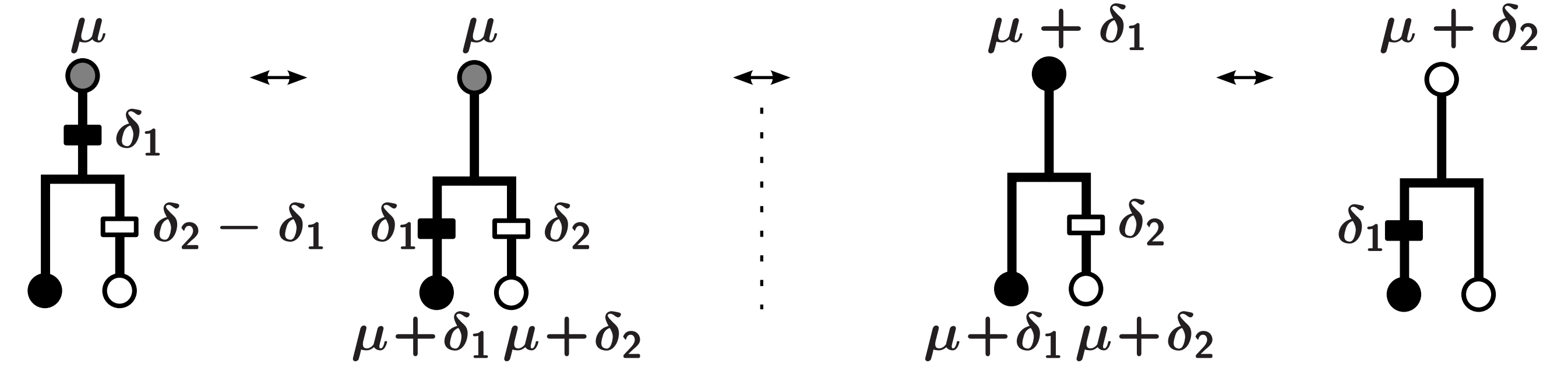
## Chelonia Dataset

[5]



**Photo Credits:** (top) "Dudhwalive chitra" by Krishna Kumar Mishra — Own work. Licensed under CC BY 3.0 via Wikimedia Commons; (middle) "Hawaii turtle 2" by Brocken Inaglor. Licensed under CC BY-SA 3.0 via Wikimedia Commons; (bottom) "Lonesome George in profile" by Mike Weston - Flickr: Lonesome George 2. Licensed under CC BY 2.0 via Wikimedia Commons.

## Identifiability Issues

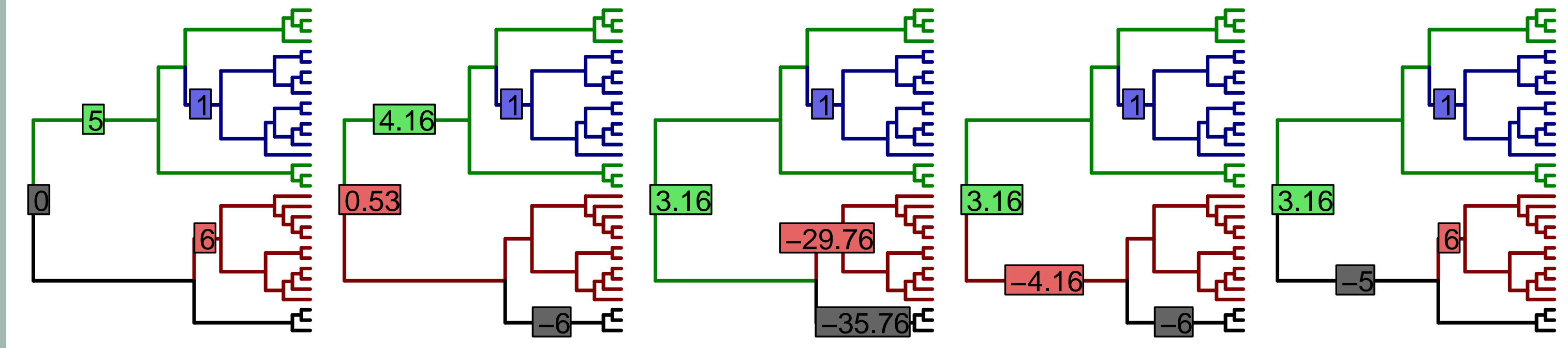


All these shifts allocations give the same tips distribution (under a BM). The two on the right are parsimonious, they are said to be *equivalent*. We discard the two on the left as over-parametrized.

## Parsimony and Equivalence Classes

Find one solution: Existing Dynamic Programming algorithms (Fitch, Sankoff) [3].

Enumerate on equivalent class: New recursive algorithm (implemented in R).



These five shifts allocations are equivalent: they are parsimonious and they produce the same tips distribution (under an OU).

## Number of Models with K shifts

No Homoplasy: 1 shift = 1 new color.

Proposition:  $K$  shifts  $\iff K + 1$  colors.

$$\mathcal{S}_K^{\text{PI}} = \{\text{Parsimonious allocations of } K \text{ shifts}\} / \sim$$

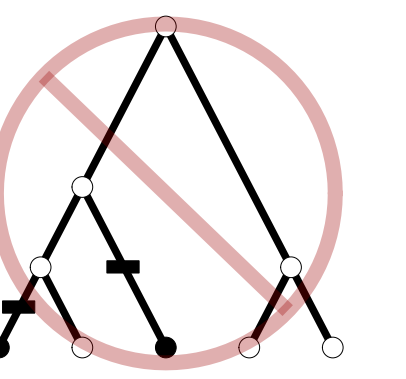
$$\mathcal{S}_K^{\text{PI}} \simeq \{\text{Coloring of tips in } K + 1 \text{ colors}\}$$

Proposition:

$$\left| \mathcal{S}_K^{\text{PI}} \right| \leq \binom{m+n-1}{K}$$

$$\left| \mathcal{S}_K^{\text{PI}} \right| \text{ depends on the topology of the tree.}$$

$$\text{► For a binary tree: } \left| \mathcal{S}_K^{\text{PI}} \right| = \binom{2n-2-K}{K}.$$



## Model Selection on K (alpha known)

[1]

Under our setting:

$$Y = R\Delta + \gamma E \quad \text{with} \quad E \sim \mathcal{N}(0, V) \quad \text{and} \quad \mathcal{S} = \{S_\eta, \eta \in \mathcal{M}\}, \quad \mathcal{M} = \bigcup_{K \geq 0} \mathcal{S}_K^{\text{PI}}$$

Define the following penalty:

$$\text{pen}(K) = A \frac{n-K-1}{n-K-2} \text{EDkhi}[K+2, n-K-2, e^{-L_K}], \quad L_K = \log |\mathcal{S}_K^{\text{PI}}| + 2 \log(K+2)$$

$$\text{and the estimator: } \hat{\eta} = \underset{\eta \in \mathcal{M}}{\text{argmin}} \|Y - \hat{s}_\eta\|_V^2 \left( 1 + \frac{\text{pen}(K_\eta)}{n - K_\eta - 1} \right)$$

Under some reasonable technical hypothesis, we get the non-asymptotic bound:

$$\mathbb{E} \left[ \frac{\|s - \hat{s}_\eta\|_V^2}{\gamma^2} \right] \leq C(A, \kappa) \left[ \inf_{\eta \in \mathcal{M}} \left\{ \frac{\|s - s_\eta\|_V^2}{\gamma^2} + D_\eta(3 + \log(n)) \right\} + 1 + \log(n) \right]$$

## Conclusion and Perspectives

- We developed a general statistical framework for trait evolution models with unconstrained shifts on ultrametric trees.
- R codes available on GitHub: <https://github.com/pbastide/Phylogenetic-EM>
- Perspectives:
  - Handle multivariate (correlated) traits.
  - Deal with uncertainty (tree, data).
  - Use fossil records (non-ultrametric tree).

## References

- [1] Y. Baraud, C. Giraud, and S. Huet. Gaussian model selection with an unknown variance. *Annals of Statistics*, 37(2):630–672, Apr. 2009. doi: 10.1214/07-AOS573.
- [2] M. A. Butler and A. A. King. Phylogenetic Comparative Analysis: A Modeling Approach for Adaptive Evolution. *The American Naturalist*, 164(6):pp. 683–695, 2004. ISSN 00030147.
- [3] J. Felsenstein. *Inferring Phylogenies*. Sinauer Associates, Sunderland, USA, 2004.
- [4] T. F. Hansen. Stabilizing selection and the comparative analysis of adaptation. *Evolution*, 51(5): 1341–1351, oct 1997.
- [5] A. L. Jaffe, G. J. Slater, and M. E. Alfaro. The evolution of island gigantism and body size variation in tortoises and turtles. *Biology letters*, 2011.