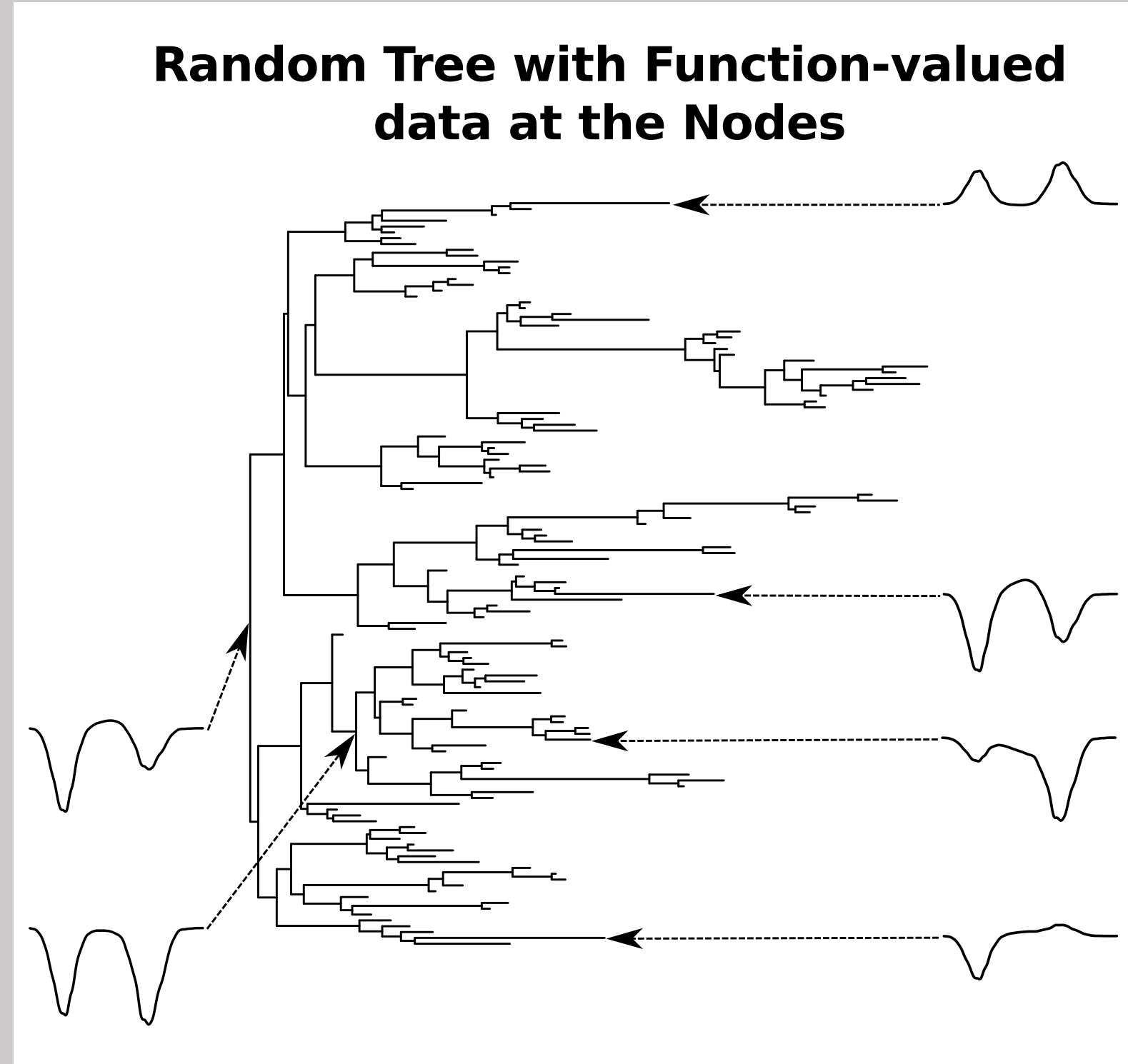


## Introduction

- ▶ Ancestral Inference is an integral part for deducing the evolutionary dynamics behind a phylogeny
  - ▷ relates to how fast a trait evolves
  - ▷ has strict physiological characteristics
  - ▷ most modelling approaches ignore curve/functional nature
- ▶ Phylogeny-specific issues
  - ▷ information is typically only available for extant organisms
  - ▷ a phylogeny describes a complex pattern of non-independence
- ▶ Current proposal
  - ▷ characterizes the trait we are interested as a curve; the realization of a stochastic gaussian process
  - ▷ identifies variations due to non-phylogenetic effects

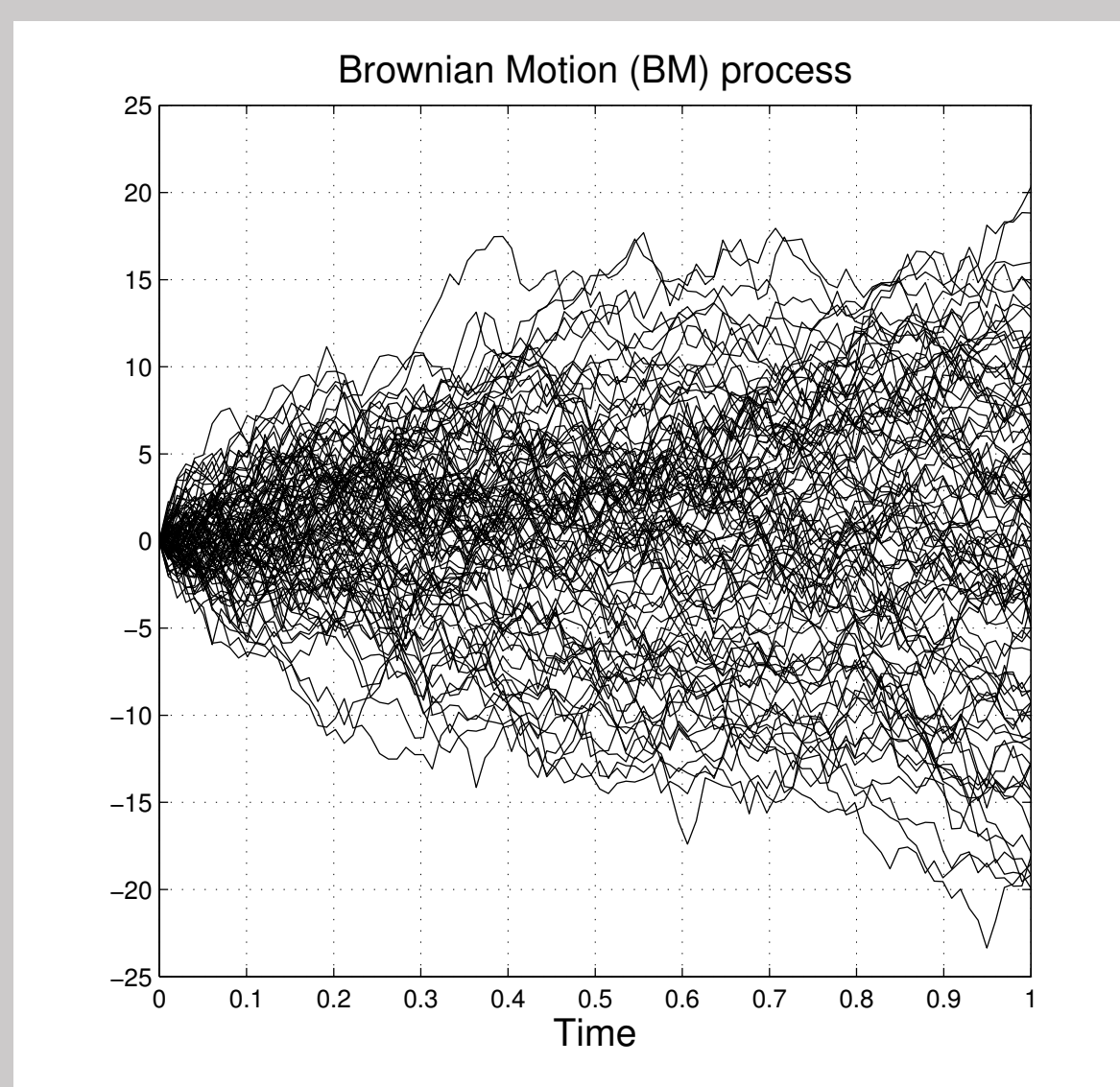


## Stochastic Processes in Evolutionary Biology

- ▶ Causes of Evolution
  - ▷ Natural Selection
  - ▷ Drift
  - ▷ Founder Effects
- ▶ By far, the most popular model evolution is that of Brownian Motion.
- ▶ Fundamental Limitation: Unaccountability for Evolution.
- ▶ Obvious Solution : Random Walk with a Deterministic Drift

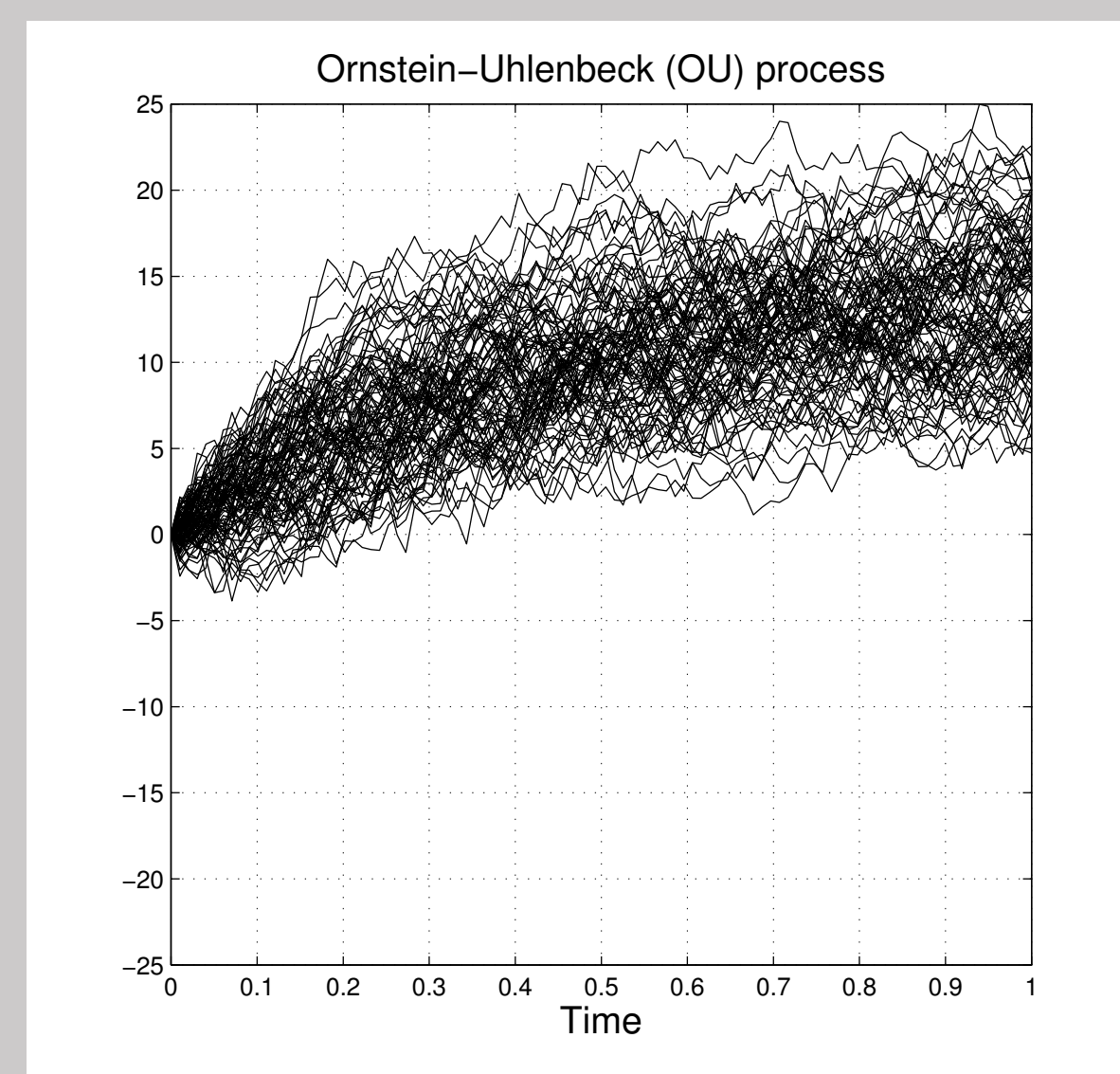
### Brownian Motion

- ▶ Markovian
- ▶  $d\mathbf{X}(t) = \sigma d\mathbf{B}(t)$
- ▶  $\sigma$ : intensity of random fluctuations in evolution
- ▶  $B(s + dt) - B(s) \sim N(0, \sigma^2 dt)$
- ▶  $\text{Cov}(B(s), B(t)) = \sigma^2 \min(s, t)$



### Ornstein-Uhlenbeck Process

- ▶ Markovian
- ▶  $d\mathbf{X}(t) = \alpha[\theta - \mathbf{X}(t)]dt + \sigma d\mathbf{B}(t)$
- ▶  $\alpha$ : strength of selection (if  $\alpha = 0 \rightarrow$  BM)
- ▶  $\theta$ : Optimum trait value
- ▶  $\text{Cov}(X(s), X(t)) = \sigma^2 \exp(-|s - t|\alpha)$

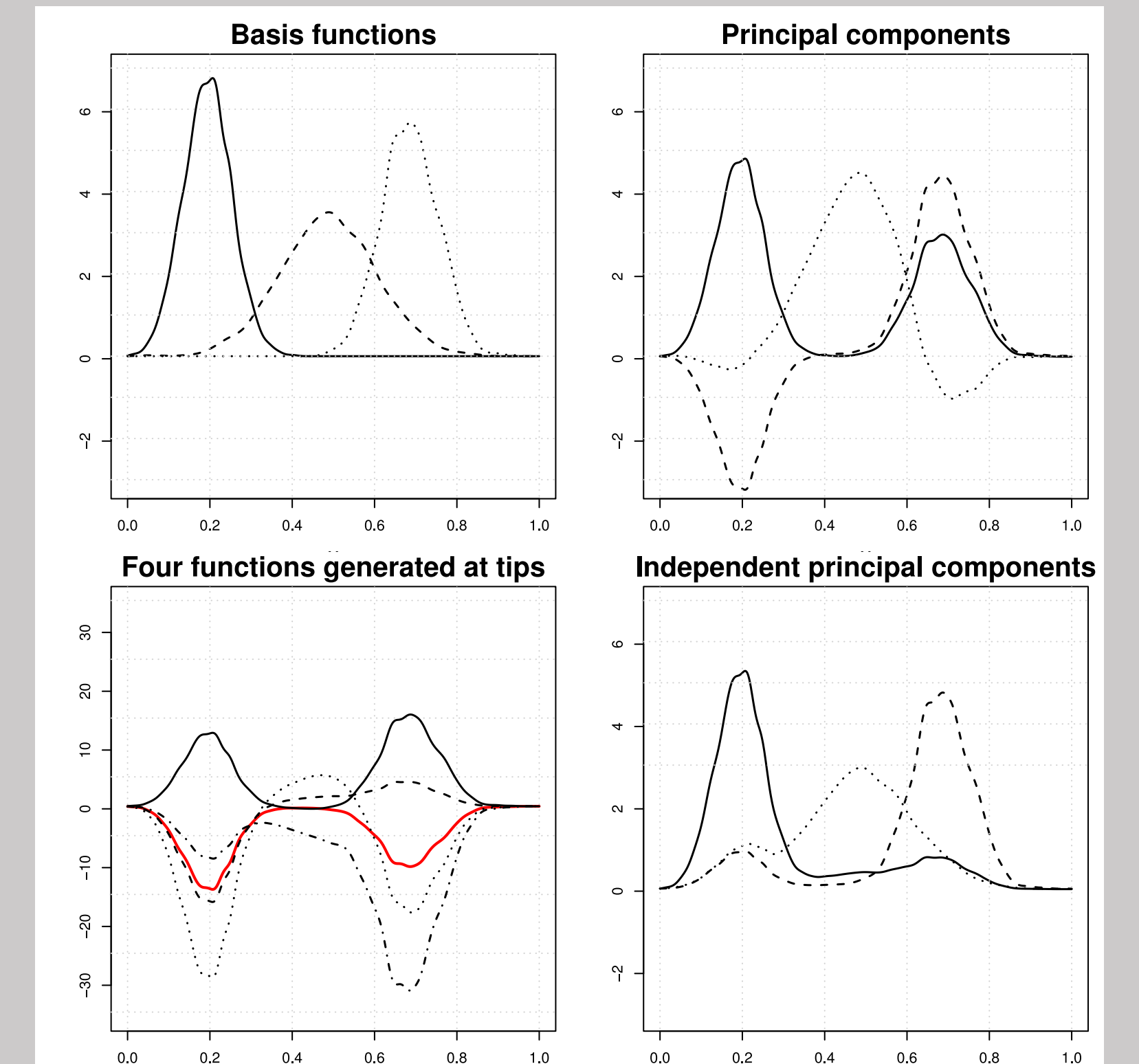


## Independent Principal Component Analysis (IPCA)

- ▶ Given a Gaussian process  $\mathbf{Y}(t), t \in [0, 1]$ , sample curves assumed to have
  - ▷  $E[\mathbf{Y}(t)] = \mu(t)$
  - ▷  $\text{Cov}[\mathbf{Y}(s), \mathbf{Y}(t)] = \mathbf{C}(s, t)$
- ▶ Mercer's Theorem for symmetric  $\mathbf{C}(s, t) = \sum_{n=1}^{\infty} \lambda_n \phi_n(s) \phi_n(t)$ 
  - ▷ order eigenvalues  $\lambda_1 \geq \lambda_2 \geq \dots \geq 0$
  - ▷ corresponding eigenfunctions  $\phi_n$ 's
- ▶  $\phi_n$ 's are by definition orthogonal; *unrealistic assumption* for real data
- ▶ Blind Source Separation
- ▶  $\Phi = \mathbf{A}\mathbf{S}$   $\mathbf{S}$ : statistically independent source signals
- ▶  $\mathbf{A}_{\text{opt}}^{-1} = \text{argmin}_{\mathbf{A}^{-1}} |\mathbf{C}_{\text{abcd}}^{(s)}|$ ,
- ▶ cumulant tensors  $\mathbf{C}_{\text{abcd}}^{(s)}$  of the output data  $\mathbf{s}_i$ .
  - ▷  $\mathbf{C}_{ijk}^{(s)} = \langle \mathbf{s}_i, \mathbf{s}_j, \mathbf{s}_k \rangle$
  - ▷  $\mathbf{C}_{ijkl}^{(s)} = \langle \mathbf{s}_i, \mathbf{s}_j, \mathbf{s}_k, \mathbf{s}_l \rangle - \langle \mathbf{s}_i, \mathbf{s}_j \rangle \langle \mathbf{s}_k, \mathbf{s}_l \rangle - \langle \mathbf{s}_i, \mathbf{s}_k \rangle \langle \mathbf{s}_j, \mathbf{s}_l \rangle - \langle \mathbf{s}_i, \mathbf{s}_l \rangle \langle \mathbf{s}_j, \mathbf{s}_k \rangle$

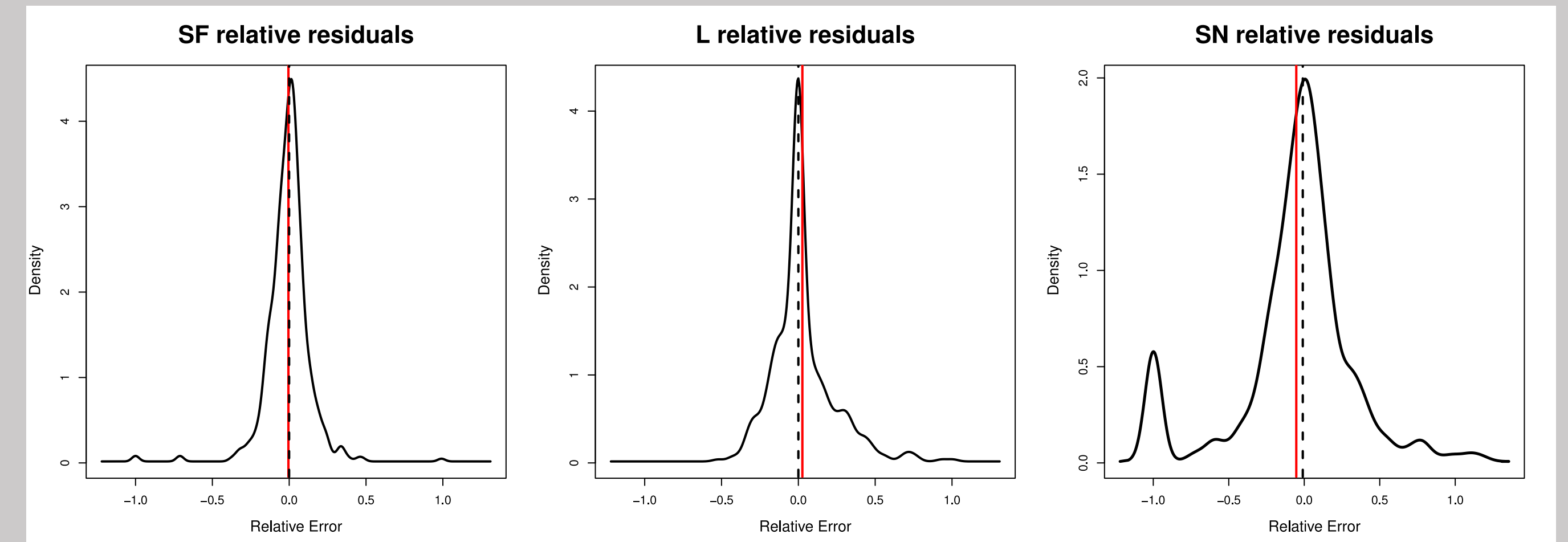
## Insights from the IPCA Analysis

- ▶ IPCA  $\mathbf{W}_{i,n}$  scores as:
  - ▷  $\hat{\mathbf{W}}_{i,n} = \sum_{k=1}^m \{ \mathbf{Y}_i(t_{i,k}) - \hat{\mathbf{m}}(t_{i,k}) \} \hat{\mathbf{s}}_n(t_{i,k}) \Delta_{i,k}$
- ▶ IPCA is a *phylogeny-agnostic*
- ▶ Automatic dimensionality determination PCA methodology
- ▶ Exclusion of outliers

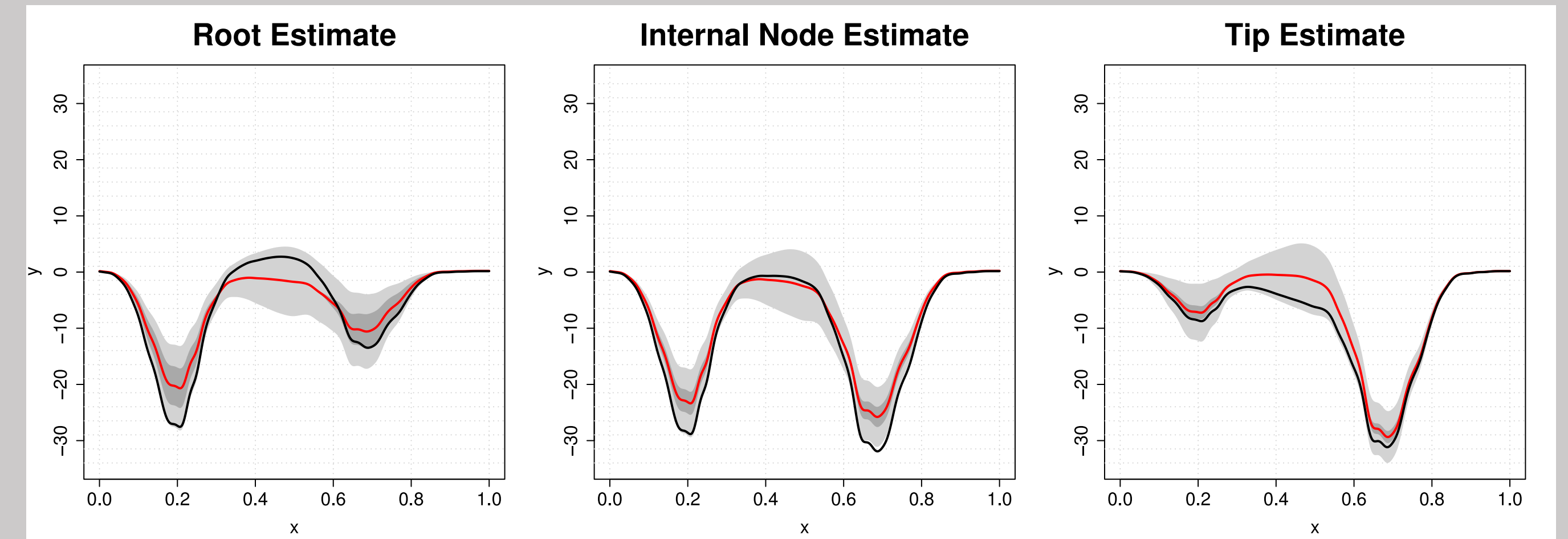


## Phylogenetic Gaussian Process Regression

- ▶ **Gaussian Process** is defined as a probability distribution over functions  $\mathbf{Y}(\mathbf{x})$  such that the set of values of  $\mathbf{Y}(\mathbf{x})$  evaluated at an arbitrary set of points  $\mathbf{x}_1, \dots, \mathbf{x}_N$  jointly have a Gaussian distribution.
- ▶ GPs are completely specified by the second-order statistics
- ▶ Assumptions
  - ▷ Conditional on their common ancestor any two trait are statistically independent. (Markov)
  - ▷ The statistical relationship between a node and any of its descendants is independent of the tree topology.
- ▶ Given trait  $\mathbf{f}(\mathbf{L})$  on a finite set of co-ordinates  $\mathbf{L}$  where  $\mathbf{K}(\mathbf{L}, \mathbf{L}, \theta)$  is the matrix of covariances of pairs  $(l_i, l_j)$  with hyperparameters  $\theta$  then :  $\mathbf{f}(\mathbf{L}) \sim \mathbf{N}(\mathbf{0}, \mathbf{K}(\mathbf{L}, \mathbf{L}, \theta))$ .
- ▶  $\mathbf{K}(l_i, l_j) = s_f^2 \exp(-|l_i - l_j|/\lambda) + s_n^2 \delta_{i,j}$ 
  - ▷  $\mathbf{K}(l_i, l_j)$ : Covariance between known points in the phylogeny
  - ▷  $s_f^2$ : intensity of random fluctuations in evolution due to balance between the restraining forces / amplitude of function variation
  - ▷  $\lambda$ : phylogenetic horizon / characteristic length scale
  - ▷  $s_n^2$ : interspecies variation, unaccountable from relations conveyed by the phylogeny / noise
- ▶ Maximizes the GP LogLikelihood:  $\log p(\mathbf{f}(\mathbf{L})|\theta) = -\frac{1}{2} \mathbf{f}(\mathbf{L})^T \mathbf{K}(\mathbf{L}, \mathbf{L}, \theta) \mathbf{f}(\mathbf{L}) - \frac{1}{2} \log |\mathbf{K}(\mathbf{L}, \mathbf{L}, \theta)| - \frac{|\mathbf{L}|}{2} \log(2\pi)$



- ▶  $\mathbf{f}(\mathbf{A})|\mathbf{f}(\mathbf{L}) \sim \mathbf{N}(\mathbf{K}(\mathbf{A}, \mathbf{L})\mathbf{K}(\mathbf{L}, \mathbf{L})^{-1}\mathbf{f}(\mathbf{L}), \mathbf{K}(\mathbf{A}) - \mathbf{K}(\mathbf{A}, \mathbf{L})\mathbf{K}(\mathbf{L}, \mathbf{L})\mathbf{K}(\mathbf{A}, \mathbf{L})^T)$



## Illustrative Bibliography

- ▶ P.Z. Hadjipantelis, N.S. Jones, J. Moriarty, D. Springate and C.G. Knight, 2012. *Ancestral Inference from Functional Data: Statistical Methods and Numerical Examples*
- ▶ N.S. Jones and J. Moriarty, 2012. *Evolutionary Inference for Function-valued Traits: Gaussian Process Regression on Phylogenies*