



HAL
open science

Fast and accurate branch length estimation for phylogenomic trees: ERaBLE (Evolutionary Rates and Branch Length Estimation)

Manuel Binet, Olivier Gascuel, Celine Scornavacca, Emmanuel J.P. Douzery, Fabio Pardi

► **To cite this version:**

Manuel Binet, Olivier Gascuel, Celine Scornavacca, Emmanuel J.P. Douzery, Fabio Pardi. Fast and accurate branch length estimation for phylogenomic trees: ERaBLE (Evolutionary Rates and Branch Length Estimation). Rencontres ALPHY - Génomique Evolutive, Bioinformatique, Alignement et Phylogénie, Mar 2015, Montpellier, France. lirmm-01237447

HAL Id: lirmm-01237447

<https://hal-lirmm.ccsd.cnrs.fr/lirmm-01237447>

Submitted on 3 Dec 2015

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Fast and accurate branch length estimation for phylogenomic trees: ERaBLE (Evolutionary Rates and Branch Length Estimation).

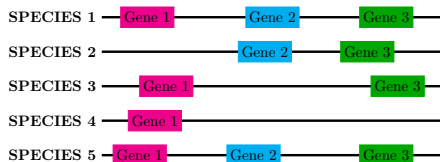
Manuel BINET, Olivier GASCUEL, Celine SCORNAVACCA, Emmanuel J.P.
DOUZERY and Fabio PARDI

March 10, 2015





Important goal: build the species tree.



New challenges:

- Use the maximum of available data.
- Account for gene rate heterogeneity.
- Deal with conflicts in gene trees topologies.

...

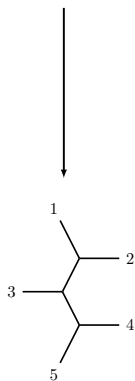
ERaBLE (Evolutionary Rates and Branch Lengths estimation).

Categories of phylogenomics methods

Based on concatenated alignment

Concatenate

1 AAGTCATACCAGCATGAC
 2 ??????ACTCCCCAGGAG
 3 AGGACC???????AAGAG
 4 GACAGA????????????
 5 GAAACCCTCTCTAAGAC



Supertree

Gene 1

1 AAGTCA
 3 AGGACC
 4 GACAGA
 5 GAAACC



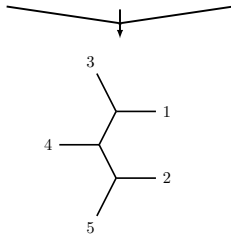
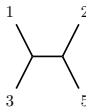
Gene 2

1 TACCAGC
 2 ACTCCCC
 5 ACTCTCT



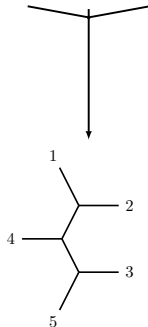
Gene 3

1 ATGAC
 2 AGGAG
 3 AAGAG
 5 AAGAC



Distance-based

δ_1 δ_2 δ_3

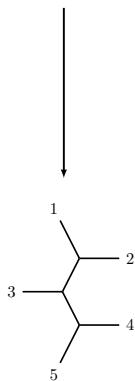
$$\begin{bmatrix} \delta_{13} \\ \delta_{14} \\ \delta_{15} \\ \delta_{34} \\ \delta_{35} \\ \delta_{45} \end{bmatrix} \begin{bmatrix} \delta_{12} \\ \delta_{15} \\ \delta_{25} \end{bmatrix} \begin{bmatrix} \delta_{12} \\ \delta_{13} \\ \delta_{15} \\ \delta_{23} \\ \delta_{25} \\ \delta_{35} \end{bmatrix}$$


Categories of phylogenomics methods

Based on concatenated alignment

Concatenate

1 AAGTCATACCAGCATGAC
 2 ??????ACTCCCCAGGAG
 3 AGGACC??????AAGAG
 4 GACAGA????????????
 5 GAAACCCTCTCTAAGAC



Branch lengths:
 yes but computationally heavy

Supertree

Gene 1

1 AAGTCA
 3 AGGACC
 4 GACAGA
 5 GAAACC



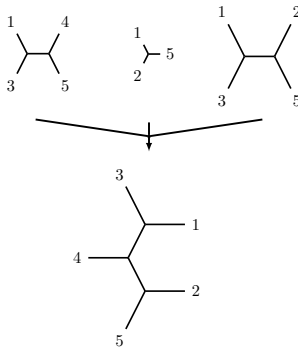
Gene 2

1 TACCAGC
 2 ACTCCCC
 5 ACTCTCT



Gene 3

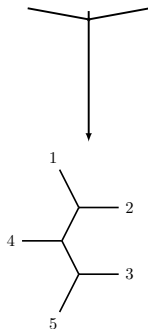
1 ATGAC
 2 AGGAG
 3 AAGAG
 5 AAGAC



usually no

Distance-based

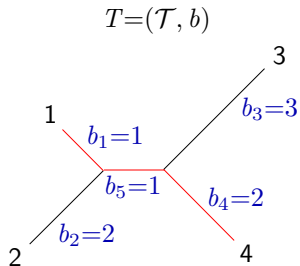
δ_1 δ_2 δ_3

$$\begin{bmatrix} \delta_{13} \\ \delta_{14} \\ \delta_{15} \\ \delta_{34} \\ \delta_{35} \\ \delta_{45} \end{bmatrix} \quad \begin{bmatrix} \delta_{12} \\ \delta_{15} \\ \delta_{25} \end{bmatrix} \quad \begin{bmatrix} \delta_{12} \\ \delta_{13} \\ \delta_{15} \\ \delta_{23} \\ \delta_{25} \\ \delta_{35} \end{bmatrix}$$


yes

Tree distances

Every tree with branch lengths can be represented with a **vector of tree distances** by computing the distance d_{ij}^T between each pair of taxa in the tree.



$$d^T = \begin{pmatrix} d_{12} = 3 \\ d_{13} = 5 \\ \mathbf{d}_{14} = \mathbf{4} \\ d_{23} = 6 \\ d_{24} = 5 \\ d_{34} = 5 \end{pmatrix}$$

Vector of tree distances

$$d_{ij}^T = \sum_{e \in P_{ij}} b_e$$

Tree distances

Every vector of tree distances is the product of the **topological matrix** of T by the branch lengths vector of T . i.e. $d^T = Ab$.

$$d^T = \begin{pmatrix} d_{12} = 3 \\ d_{13} = 5 \\ d_{14} = 4 \\ d_{23} = 6 \\ d_{24} = 5 \\ d_{34} = 5 \end{pmatrix}$$

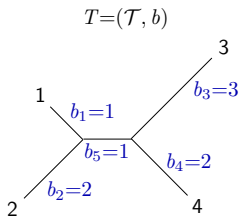
Vector of tree distances

$$A = \begin{pmatrix} 1 & 1 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 1 \\ 1 & 0 & 0 & 1 & 1 \\ 0 & 1 & 1 & 0 & 1 \\ 0 & 1 & 0 & 1 & 1 \\ 0 & 0 & 1 & 1 & 0 \end{pmatrix}$$

Topological matrix

$$b = \begin{pmatrix} b_1 = 1 \\ b_2 = 2 \\ b_3 = 3 \\ b_4 = 2 \\ b_5 = 1 \end{pmatrix}$$

Branch lengths vector



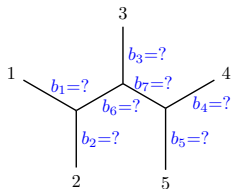
$$d^T = Ab$$

A distance-based method : the weighted least square method (WLS)

Branch lengths estimation with WLS:

$$\delta = \begin{bmatrix} \delta_{12} = 0.140 \\ \delta_{13} = 0.163 \\ \delta_{14} = 0.288 \\ \delta_{15} = 0.336 \\ \delta_{23} = 0.188 \\ \delta_{24} = 0.413 \\ \delta_{25} = 0.411 \\ \delta_{34} = 0.298 \\ \delta_{35} = 0.213 \end{bmatrix}$$

input distances
(estimated from
1 gene alignment)



input topology

$$d^T = Ab$$

$$\text{minimize } \sum_{ij} w_{ij} (\delta_{ij} - d_{ij}^T)^2$$

A distance-based method : the weighted least square method (WLS)

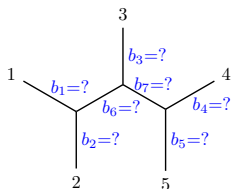
Branch lengths estimation with WLS:

$$\delta = \begin{bmatrix} \delta_{12} = 0.140 \\ \delta_{13} = 0.163 \\ \delta_{14} = 0.288 \\ \delta_{15} = 0.336 \\ \delta_{23} = 0.188 \\ \delta_{24} = 0.413 \\ \delta_{25} = 0.411 \\ \delta_{34} = 0.298 \\ \delta_{35} = 0.213 \end{bmatrix}$$

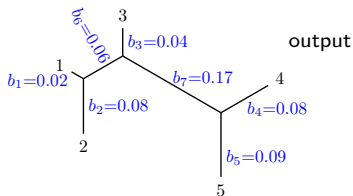
input distances
(estimated from
1 gene alignment)

$$\text{minimize } \sum_{ij} w_{ij} (\delta_{ij} - d_{ij}^T)^2$$

can be solved analytically in $\mathcal{O}(n^3)$

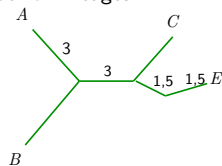
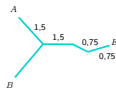
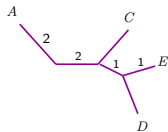


$$d^T = Ab$$



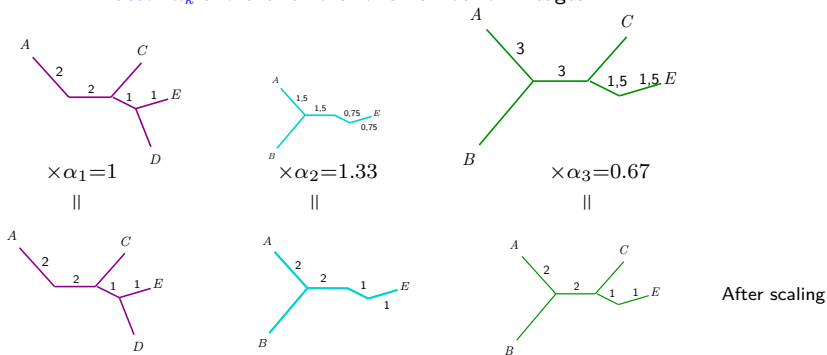
New methode ERaBLE: Evolutionary Rates and Branch Length Estimation

Hypothesis: Any gene G_k induces approximately the same tree up to a scale factor α_k and the removal of a number of lineages.



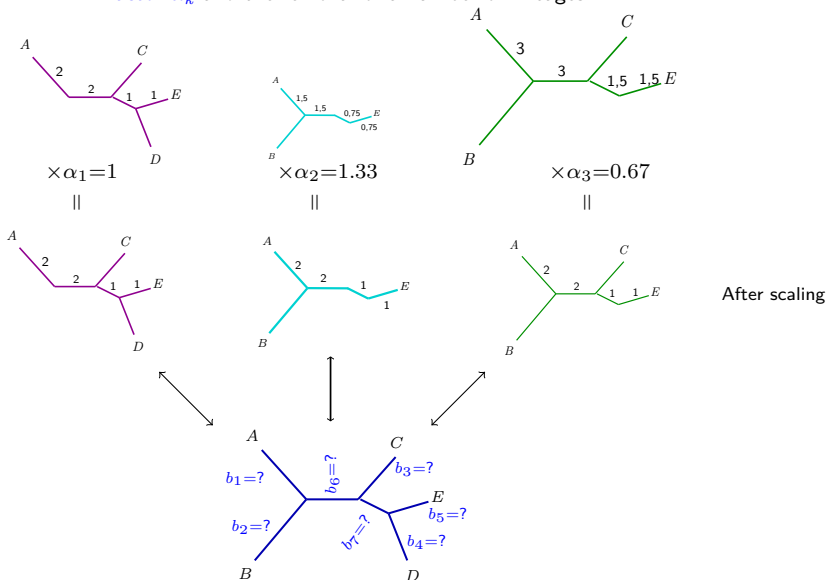
New method ERaBLE: Evolutionary Rates and Branch Length Estimation

Hypothesis: Any gene G_k induces approximately the same tree up to a **scale factor** α_k and the removal of a number of lineages.



New method ERaBLE: Evolutionary Rates and Branch Length Estimation

Hypothesis: Any gene G_k induces approximately the same tree up to a **scale factor** α_k and the removal of a number of lineages.



New methode ERaBLE: Evolutionary Rates and Branch Length Estimation

Input data: m distance vectors $\delta_1, \dots, \delta_m$ where $\delta_k = (\delta_{ij}^{(k)})$ is defined on the taxa set L_k .

A given topology \mathcal{T} defined on the taxa set $L = \bigcup_{k=1}^m L_k$.

New method ERaBLE: Evolutionary Rates and Branch Length Estimation

Input data: m distance vectors $\delta_1, \dots, \delta_m$ where $\delta_k = (\delta_{ij}^{(k)})$ is defined on the taxa set L_k .

A given topology \mathcal{T} defined on the taxa set $L = \bigcup_{k=1}^m L_k$.

Objective: find

- the branch lengths $\hat{b}_1, \dots, \hat{b}_{2n-3}$ of \mathcal{T} ,
- the scale factors $\hat{\alpha}_1, \dots, \hat{\alpha}_m$ of the m genes,

solution of the problem:

$$\left\{ \begin{array}{l} \text{minimize} \\ \sum_{k=1}^m \sum_{i,j \in L_k} w_{ij}^{(k)} (\hat{\alpha}_k \delta_{ij}^{(k)} - d_{ij}^T)^2 \end{array} \right.$$

New method ERaBLE: Evolutionary Rates and Branch Length Estimation

Input data: m distance vectors $\delta_1, \dots, \delta_m$ where $\delta_k = (\delta_{ij}^{(k)})$ is defined on the taxa set L_k .

A given topology \mathcal{T} defined on the taxa set $L = \bigcup_{k=1}^m L_k$.

Objective: find

- the branch lengths $\hat{b}_1, \dots, \hat{b}_{2n-3}$ of \mathcal{T} ,
- the scale factors $\hat{\alpha}_1, \dots, \hat{\alpha}_m$ of the m genes,

solution of the problem:

$$\left\{ \begin{array}{l} \text{minimize} \\ \text{subject to} \end{array} \right. \quad \begin{array}{l} \sum_{k=1}^m \sum_{i,j \in L_k} w_{ij}^{(k)} (\hat{\alpha}_k \delta_{ij}^{(k)} - d_{ij}^T)^2 \\ \sum_{k=1}^m Z_k \hat{\alpha}_k = 1 \end{array}$$

New method ERaBLE: Evolutionary Rates and Branch Length Estimation

Input data: m distance vectors $\delta_1, \dots, \delta_m$ where $\delta_k = (\delta_{ij}^{(k)})$ is defined on the taxa set L_k .

A given topology \mathcal{T} defined on the taxa set $L = \bigcup_{k=1}^m L_k$.

Objective: find

- the branch lengths $\hat{b}_1, \dots, \hat{b}_{2n-3}$ of \mathcal{T} ,
- the scale factors $\hat{\alpha}_1, \dots, \hat{\alpha}_m$ of the m genes,

solution of the problem:

$$\left\{ \begin{array}{l} \text{minimize} \\ \text{subject to} \end{array} \right. \quad \begin{array}{l} \sum_{k=1}^m \sum_{i,j \in L_k} w_{ij}^{(k)} (\hat{\alpha}_k \delta_{ij}^{(k)} - d_{ij}^T)^2 \\ \sum_{k=1}^m Z_k \hat{\alpha}_k = 1 \end{array}$$

Outputs:

$$\hat{b}_e \quad \text{et} \quad \hat{r}_k = \frac{1}{\hat{\alpha}_k}$$

Resolution of the optimization problem

The minimization problem can be solved with the method of Lagrange multipliers and leads to the linear system in $\mathcal{O}(n + m)$ equations and unknowns:

$$\begin{pmatrix}
 \delta_1^t W_1 \delta_1 & 0 & \cdots & 0 & [-\delta_1^t W_1 A_1] & 1 \\
 0 & \delta_2^t W_2 \delta_2 & & \vdots & [-\delta_2^t W_2 A_2] & 1 \\
 \vdots & & \ddots & \vdots & \vdots & \vdots \\
 0 & \cdots & \cdots & \delta_m^t W_m \delta_m & [-\delta_m^t W_m A_m] & 1 \\
 \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\
 -A_1^t W_1 \delta_1 & -A_2^t W_2 \delta_2 & \cdots & -A_m^t W_m \delta_m & \left[\sum_{k=1}^m A_k^t W_k A_k \right] & 0 \\
 \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\
 Z_1 & Z_2 & \cdots & Z_m & 0 \cdots 0 & 0 \\
 & & & & & \lambda
 \end{pmatrix}
 \begin{pmatrix}
 \alpha_1 \\
 \alpha_2 \\
 \vdots \\
 \alpha_m \\
 b_1 \\
 b_2 \\
 \vdots \\
 b_{2n-3} \\
 \lambda
 \end{pmatrix}
 =
 \begin{pmatrix}
 0 \\
 \vdots \\
 \vdots \\
 \vdots \\
 \vdots \\
 \vdots \\
 \vdots \\
 0 \\
 1
 \end{pmatrix}$$

- Filling the matrix and standard system solving in $\mathcal{O}((n + m)^3 + mn^4)$

- Solving in $\mathcal{O}(n^3 + mn^2)$ with our algorithms:

- Filling the matrix in $\mathcal{O}(mn^2)$ instead of $\mathcal{O}(mn^4)$
- Block-solving system in $\mathcal{O}(n^3)$ instead of $\mathcal{O}(n + m)^3$

$m = \#$ genes
 $n = \#$ taxa

Resolution of the optimization problem

The minimization problem can be solved with the method of Lagrange multipliers and leads to the linear system in $\mathcal{O}(n + m)$ equations and unknowns:

$$\begin{pmatrix}
 \delta_1^t W_1 \delta_1 & 0 & \cdots & 0 & [-\delta_1^t W_1 A_1] & 1 \\
 0 & \delta_2^t W_2 \delta_2 & & \vdots & [-\delta_2^t W_2 A_2] & 1 \\
 \vdots & & \ddots & \vdots & \vdots & \vdots \\
 0 & \cdots & \cdots & \delta_m^t W_m \delta_m & [-\delta_m^t W_m A_m] & 1 \\
 \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\
 -A_1^t W_1 \delta_1 & -A_2^t W_2 \delta_2 & \cdots & -A_m^t W_m \delta_m & \left[\sum_{k=1}^m A_k^t W_k A_k \right] & 0 \\
 \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\
 Z_1 & Z_2 & \cdots & Z_m & 0 \cdots 0 & 0 \\
 & & & & & \lambda
 \end{pmatrix}
 \begin{pmatrix}
 \alpha_1 \\
 \alpha_2 \\
 \vdots \\
 \alpha_m \\
 b_1 \\
 b_2 \\
 \vdots \\
 b_{2n-3} \\
 \lambda
 \end{pmatrix}
 =
 \begin{pmatrix}
 0 \\
 \vdots \\
 \vdots \\
 \vdots \\
 \vdots \\
 \vdots \\
 \vdots \\
 0 \\
 1
 \end{pmatrix}$$

- Filling the matrix and standard system solving in $\mathcal{O}((n + m)^3 + mn^4)$

- Solving in $\mathcal{O}(n^3 + mn^2)$ with our algorithms:

- Filling the matrix in $\mathcal{O}(mn^2)$ instead of $\mathcal{O}(mn^4)$
- Block-solving system in $\mathcal{O}(n^3)$ instead of $\mathcal{O}(n + m)^3$

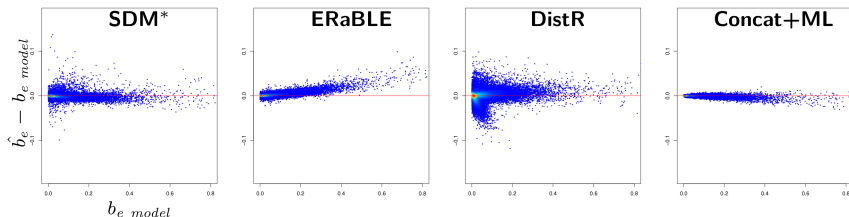
$m = \#$ genes
 $n = \#$ taxa

Complexity: ERaBLE $\mathcal{O}(n^3 + mn^2)$ vs. WLS $\mathcal{O}(n^3)$

Results on a simulated dataset

input data: $m=500$ genes defined on $n \in [4, 40]$ taxa + model topology (500 replicates).

Accuracy in the estimation of branch lengths:



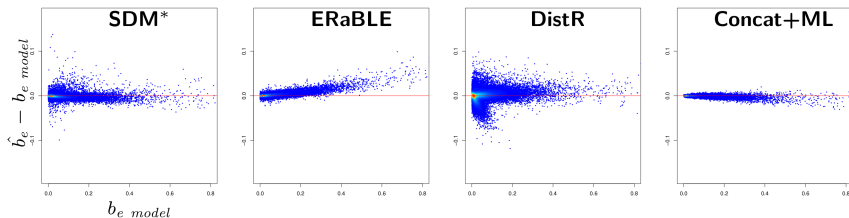
Compared methods (after adaptation):

- SDM* [Criscuolo et al. 2006]. Phylogenomic distance-based method,
- DistR [Bevan et al. 2005]. Distance-based method for the estimation of gene rates,
- Concat+ML. PhyML [Guindon et al. 2010] analysis of the concatenated alignments.

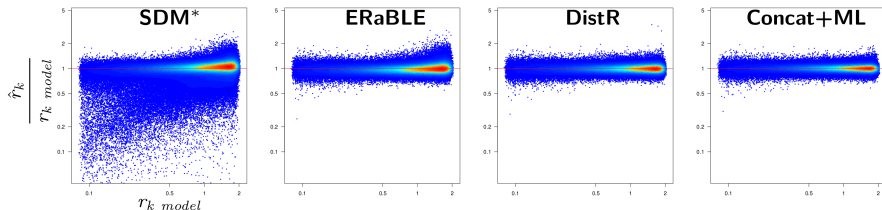
Results on a simulated dataset

input data: $m=500$ genes defined on $n \in [4, 40]$ taxa + model topology (500 replicates).

Accuracy in the estimation of branch lengths:



Accuracy in the estimation of gene rates (logarithmic scale):

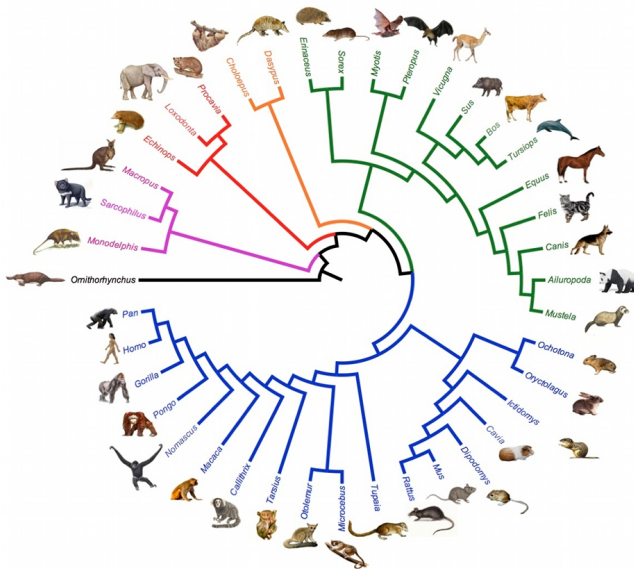


Results on the OrthoMaM dataset [Douzery et al. 2014]



input data: $m=6953$ nucleotide exon alignments over $n=4$ to 40 mammals.

input topology: the topology of the 40 mammals present in OrthoMaM:

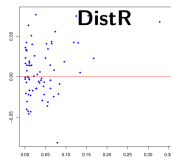
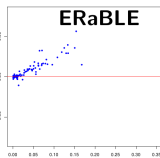
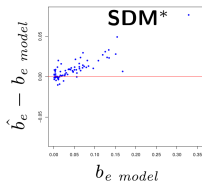


Results on the OrthoMaM dataset

	SDM*	ERaBLE	DistR	Concat+ML
Running times	8h33m*	7s*	2h8m*	41h16m
RAM	1.2 GB	221 MB	3 GB	488 GB

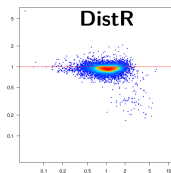
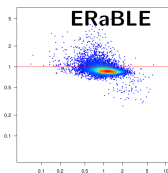
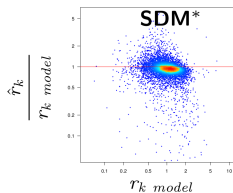
* add 2min46s for the 6953 input distances estimation

Accuracy in the estimation of branch lengths:



Concat+ML
as model

Accuracy in the estimation of gene rates (logarithmic scale):



Concat+ML
as model

Conclusion

ERaBLE gives branch lengths to phylogenomic trees (e.g. estimated with supertree methods).

ERaBLE is relatively accurate in the estimations of branch lengths and of gene rates.

ERaBLE is fast with a complexity in $\mathcal{O}(n^3 + mn^2)$ (linear in m).

Thank you for your attention.

Any questions?

Fundings:

