

# Influence Function for Robust Phylogenetic Reconstruction

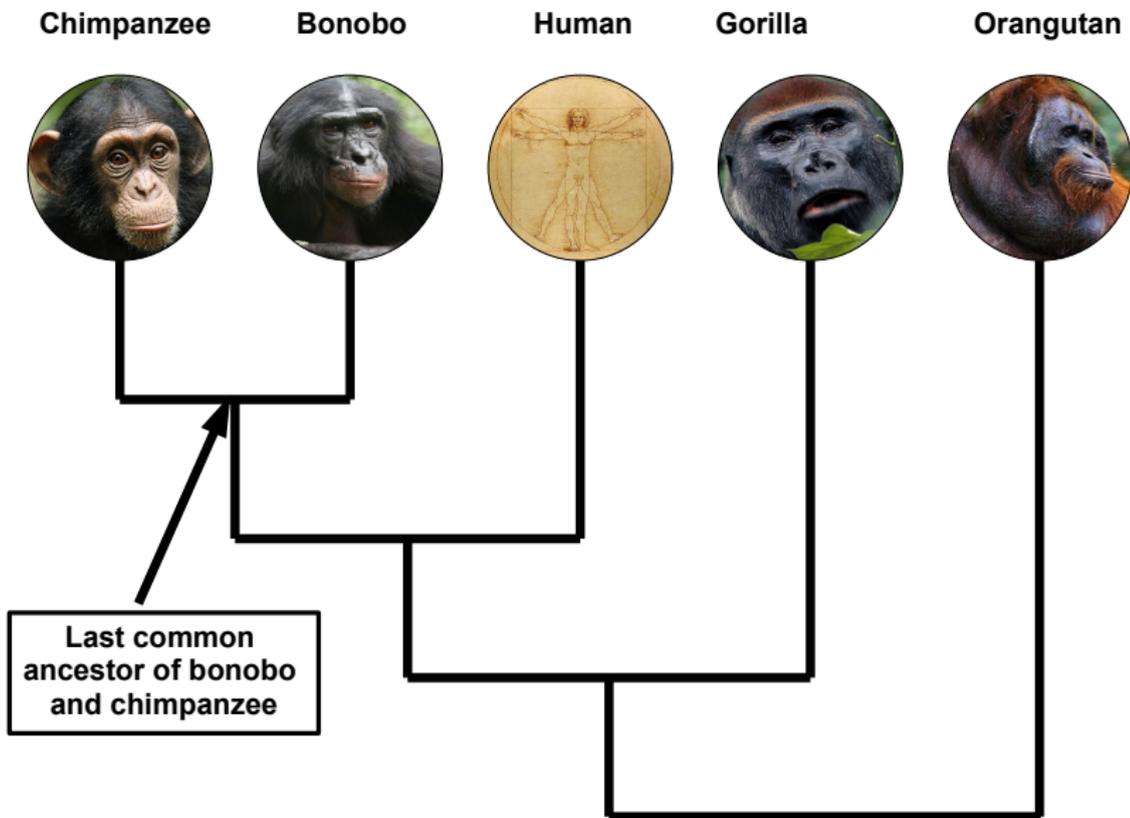
Mahendra Mariadassou

Unité MIG  
INRA Jouy-en-Josas

JOBIM 2011  
Institut Pasteur

- 1 Motivation
- 2 Outlier Sites
- 3 Taxon Influence Index

# Goal of Phylogenetic Reconstruction



# Applications in Many Domains

Including but not limited to:

## Studies of an Evolutionary Group:

- Estimate Time to Most Recent Common Ancestor (TMRCA);
- Find genes under positive/purifying selection;
- Identifying Horizontal Gene Transfer;
- Testing evolutionary hypothesis.

## Systematics:

- Reconstruct Tree of X, phylogeny of all living X;
- DNA barcoding: easily identify the species of a new organism;
- Natural way to measure biodiversity.

Most of these applications require “good” trees.

# Applications in Many Domains

Including but not limited to:

## Studies of an Evolutionary Group:

- Estimate Time to Most Recent Common Ancestor (TMRCA);
- Find genes under positive/purifying selection;
- Identifying Horizontal Gene Transfer;
- Testing evolutionary hypothesis.

## Systematics:

- Reconstruct Tree of X, phylogeny of all living X;
- DNA barcoding: easily identify the species of a new organism;
- Natural way to measure biodiversity.

Most of these applications require “good” trees.

# Applications in Many Domains

Including but not limited to:

## Studies of an Evolutionary Group:

- Estimate Time to Most Recent Common Ancestor (TMRCA);
- Find genes under positive/purifying selection;
- Identifying Horizontal Gene Transfer;
- Testing evolutionary hypothesis.

## Systematics:

- Reconstruct Tree of X, phylogeny of all living X;
- DNA barcoding: easily identify the species of a new organism;
- Natural way to measure biodiversity.

Most of these applications require “good” trees.

# Reconstruction Methods

## Three main families:

**Distance based:** Neighbor-Joining (NJ),

**Parsimony based:** Maximum Parsimony (MP),

**Likelihood based:** Maximum Likelihood (ML), Bayesian Inference (BI).

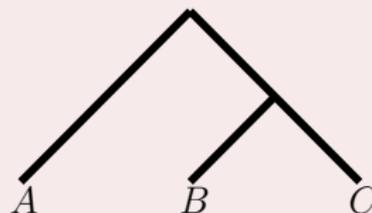
### Alignment

<b>A</b>	<i>C</i>	<i>C</i>	<i>T</i>	<i>T</i>
<b>B</b>	<i>G</i>	<i>G</i>	<i>A</i>	<i>A</i>
<b>C</b>	<i>G</i>	<i>G</i>	<i>A</i>	<i>C</i>



NJ, MP,  
ML, BI,  
...

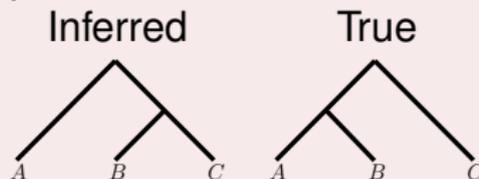
### Phylogenetic Tree



# Validating the Tree

## Inference Problems:

- Compare **inferred tree** to **true tree** to assess how good it is,



- But the true tree is not available!

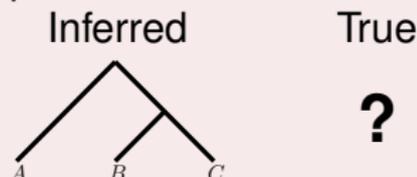
## Confidence Issue:

- How **confident** are we on the inferred tree ?
- Which **parts** of the tree are **reliable/not reliable** ?
- How **robust** is the tree to small changes in the data and **outliers** ?

# Validating the Tree

## Inference Problems:

- Compare **inferred tree** to **true tree** to assess how good it is,



- **But the true tree is not available!**

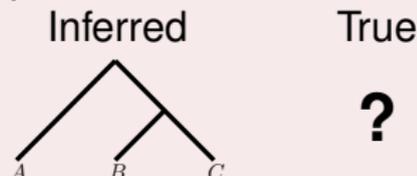
## Confidence Issue:

- How **confident** are we on the inferred tree ?
- Which **parts** of the tree are **reliable/not reliable** ?
- How **robust** is the tree to small changes in the data and **outliers** ?

# Validating the Tree

## Inference Problems:

- Compare **inferred tree** to **true tree** to assess how good it is,



- **But the true tree is not available!**

## Confidence Issue:

- How **confident** are we on the inferred tree ?
- Which **parts** of the tree are **reliable/not reliable** ?
- How **robust** is the tree to small changes in the data and **outliers** ?

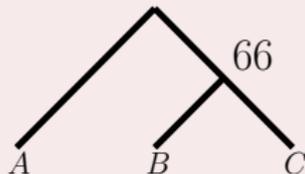
# Bootstrap Values: the Theory

## Original Dataset:

### Alignment

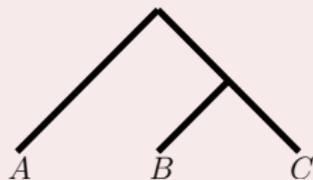
<b>A</b>	A	C	T	T
<b>B</b>	G	G	A	T
<b>C</b>	G	G	C	C

### Phylogenetic tree

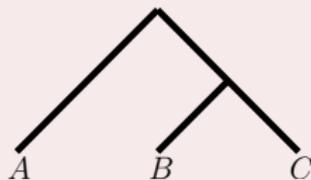


## Bootstrap Datasets:

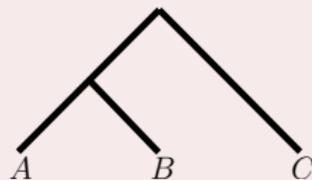
<b>A</b>	A	C	T	C
<b>B</b>	G	G	A	G
<b>C</b>	G	G	C	G



<b>A</b>	C	A	T	A
<b>B</b>	G	G	A	G
<b>C</b>	G	G	C	G



<b>A</b>	T	T	T	T
<b>B</b>	A	T	A	T
<b>C</b>	C	C	C	C



# Bootstrap Values: A Robustness Index ?

## Potential Causes for Uncertainty:

- Finite sequence lengths (sampling errors);
- Poor alignment quality (influent sites);
- Poor taxon sampling (rogue taxa);
- Model misspecification,
- ...

## Different Tools to Assess Them:

- Bootstrap: deals with **global** sources of uncertainty;
- Unable to pinpoint **local** sources of uncertainty;
- Need for other indexes to detect **outliers**

# Bootstrap Values: A Robustness Index ?

## Potential Causes for Uncertainty:

- Finite sequence lengths (sampling errors);
- Poor alignment quality (influent sites);
- Poor taxon sampling (rogue taxa);
- Model misspecification,
- ...

## Different Tools to Assess Them:

- Bootstrap: deals with global sources of uncertainty;
- Unable to pinpoint local sources of uncertainty;
- Need for other indexes to detect outliers

## Motivation: Filter Data

Identifying/filter out **outliers** corresponding to:

- Sequencing errors;
- Alignment errors;
- Presence of an atypical DNA segment;
- ...

## Procedure

- Quantify **influence** of site  $i$  by
  - Removing site  $i$  from alignment;
  - Computing new tree  $T^{-i}$  from smaller (jackknife) alignment;
  - Comparing new tree to original tree;
- Compute phylogeny from “**not too influent sites**”.

## Motivation: Filter Data

Identifying/filter out **outliers** corresponding to:

- Sequencing errors;
- Alignment errors;
- Presence of an atypical DNA segment;
- ...

## Procedure

- Quantify **influence** of site  $i$  by
  - Removing site  $i$  from alignment;
  - Computing new tree  $T^{-i}$  from smaller (jackknife) alignment;
  - Comparing new tree to original tree;
- Compute phylogeny from “**not too influent sites**”.

## Definition

- For alignment of size  $n$ , the **influence value** of site  $i$  is:

$$IF(\text{site } i) = (n - 1) \left( \frac{\text{LogLik}(T^{-i})}{n - 1} - \frac{\text{LogLik}(T)}{n} \right)$$

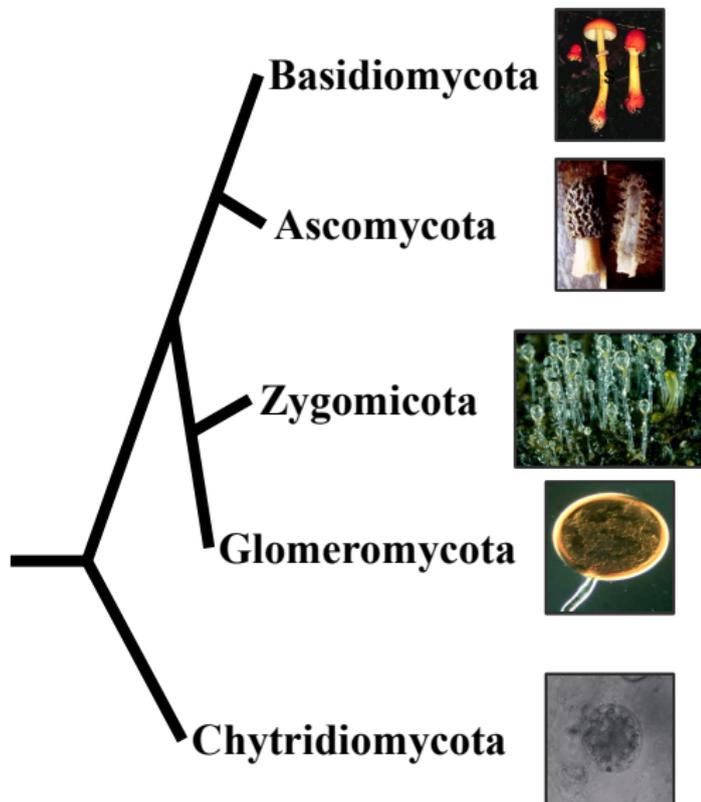
- **Difference** between support of alignments for their ML trees;
- Negative / Positive value: **enhanced** / **weakened** support (when adding the site);
- Expect most sites with **small negative** values and a few with **large positive** values.

## Strategy towards greater stability

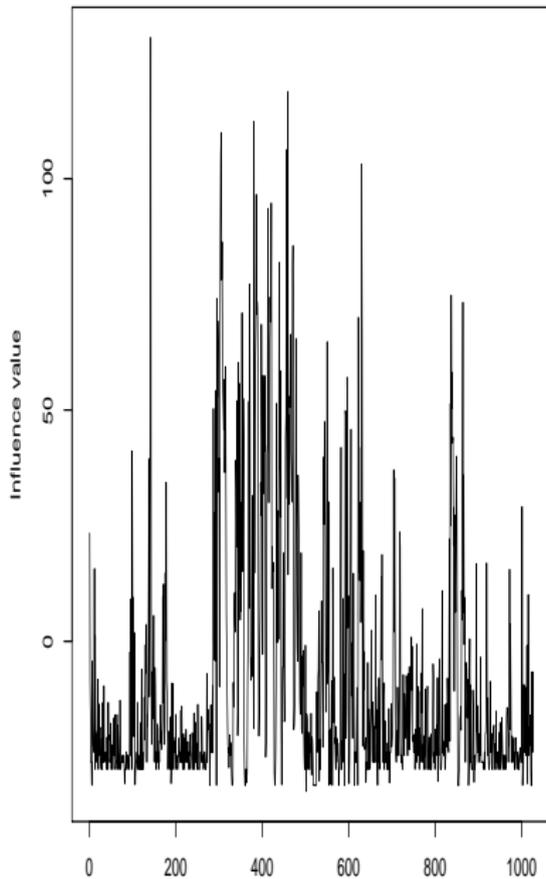
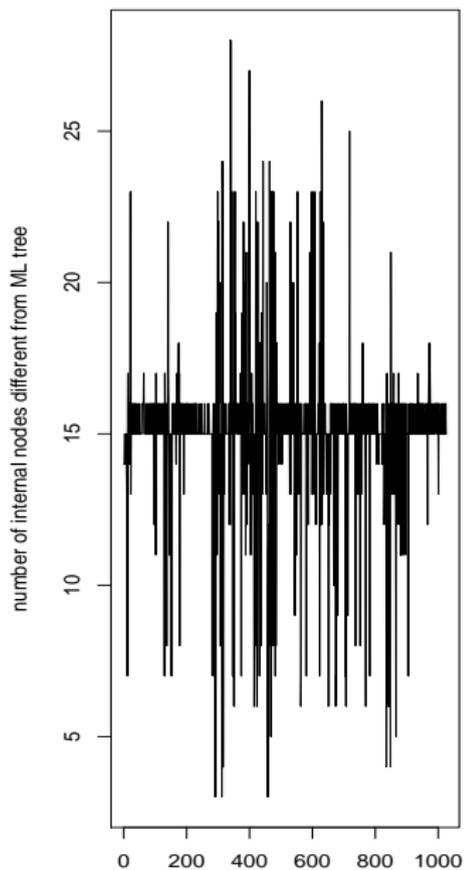
- Focus on **outliers**: sites with  $IF(\text{site } i) > 0$ ;
- Rank them in decreasing  $IF(\text{site } i)$ ;
- Remove them one at the time until a stable tree is found.

# Data: Zygomycetes & Chytridiomycetes

- Lower mushrooms;
- Biology: unknown!
- Enough signal to resolve the topology;
- 1026 sites, 158 taxa, GTR model.

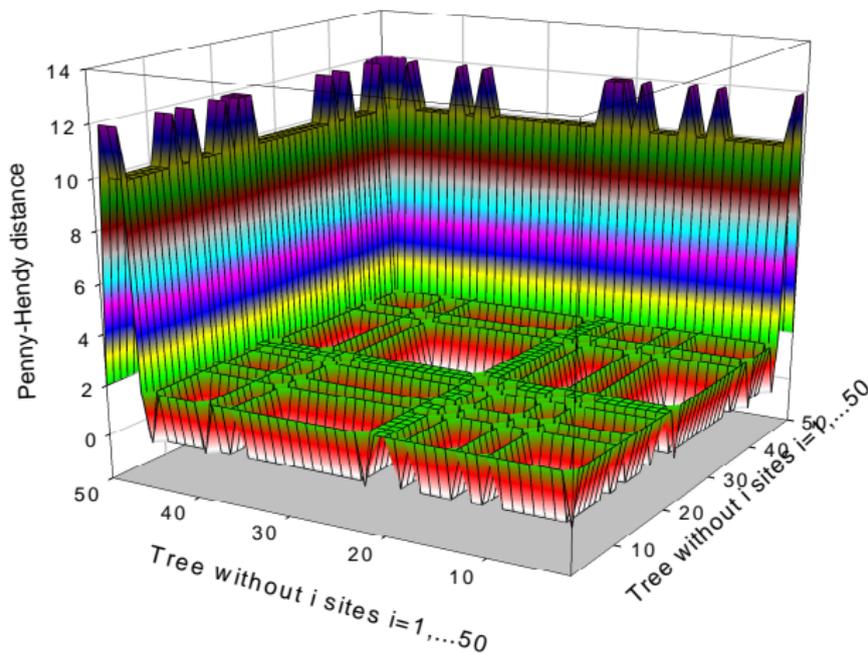


# Information About Sites



# Distance Between Trees

Distance between trees



$$d(T_0, T_i) \approx 18 \quad \text{and} \quad d(T_i, T_j) \leq 2 \quad \text{for} \quad i, j = 1..45$$

- Strongest Outlier (position 142): Highly variable site located on a med loop (5nt) located on a conserved hairpin;
- Removing most influent sites leads to Increased bootstrap values and loss of 20% of inner nodes;
- Confirms monophyly of phyla Glomeromycota
- Reinforces polyphyletic status of phyla Chytridiomycota and Zygomycota

## Motivation: Filter Data

- Study the robustness of the tree with respect to the species
- Identify **rogue** taxa.

## Procedure

- Quantify **influence** of taxon  $i$  by:
  - Removing site  $i$  from alignment;
  - Computing new tree  $T^{-i}$  from smaller (jackknife) alignment;
  - Comparing new tree to original tree;
- Compute phylogeny from “**not too influent sites**”.

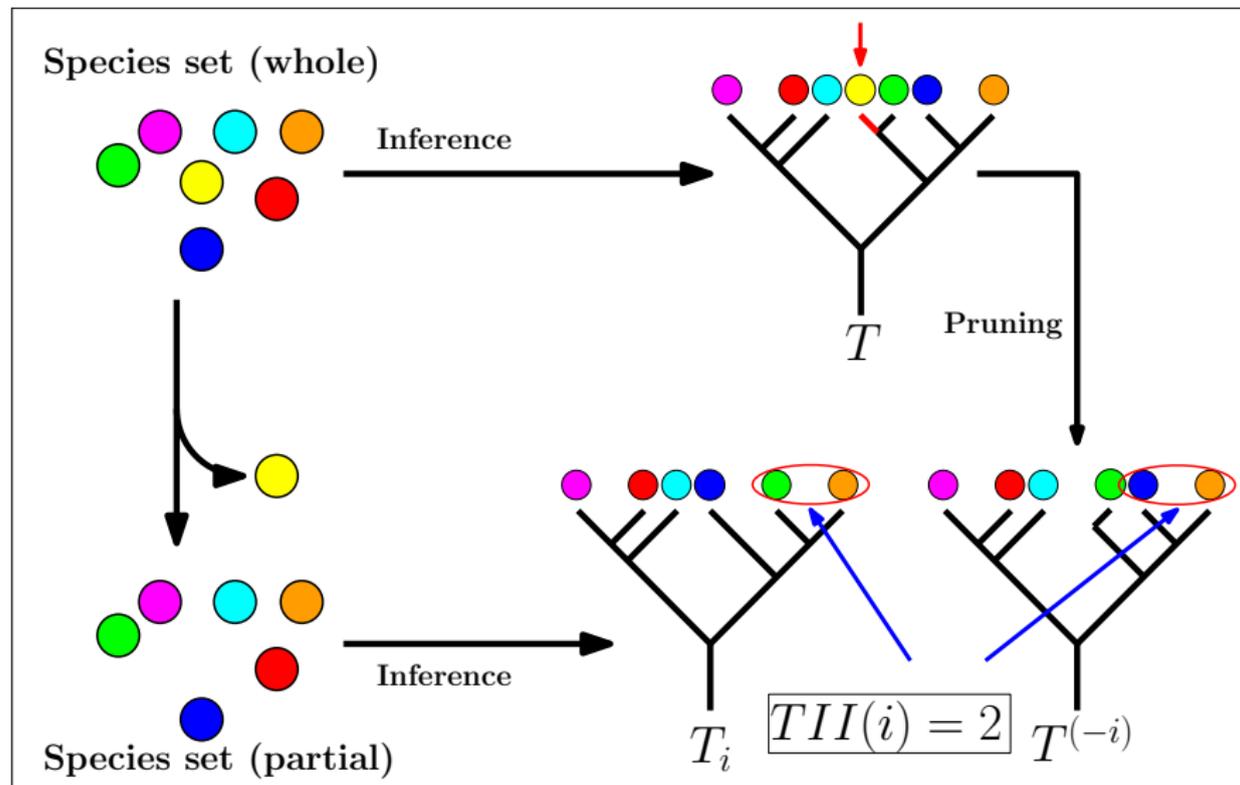
## Motivation: Filter Data

- Study the robustness of the tree with respect to the species
- Identify **rogue** taxa.

## Procedure

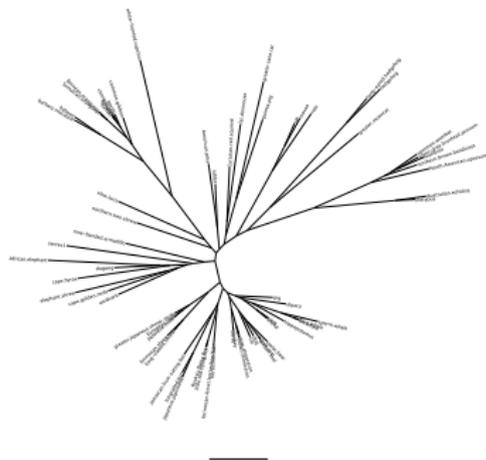
- Quantify **influence** of taxon  $i$  by:
  - Removing site  $i$  from alignment;
  - Computing new tree  $T^{-i}$  from smaller (jackknife) alignment;
  - Comparing new tree to original tree;
- Compute phylogeny from “**not too influent sites**”.

# Method

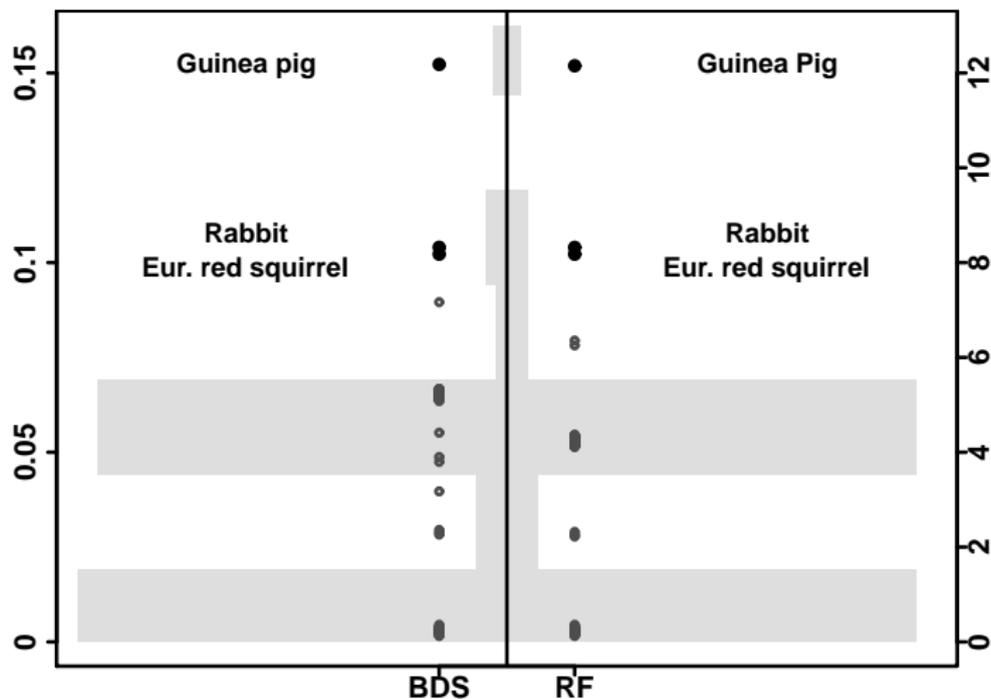


# Data: Placental Mammal Phylogeny

- Mitochondrial genome of 68 mammals;
- Amino Acids sequences;
- Sequences are 3658 AA long;
- MtMam + I +  $\Gamma$ 4 model;
- Phylogeny published in Nikaido *et al.* in 2003.

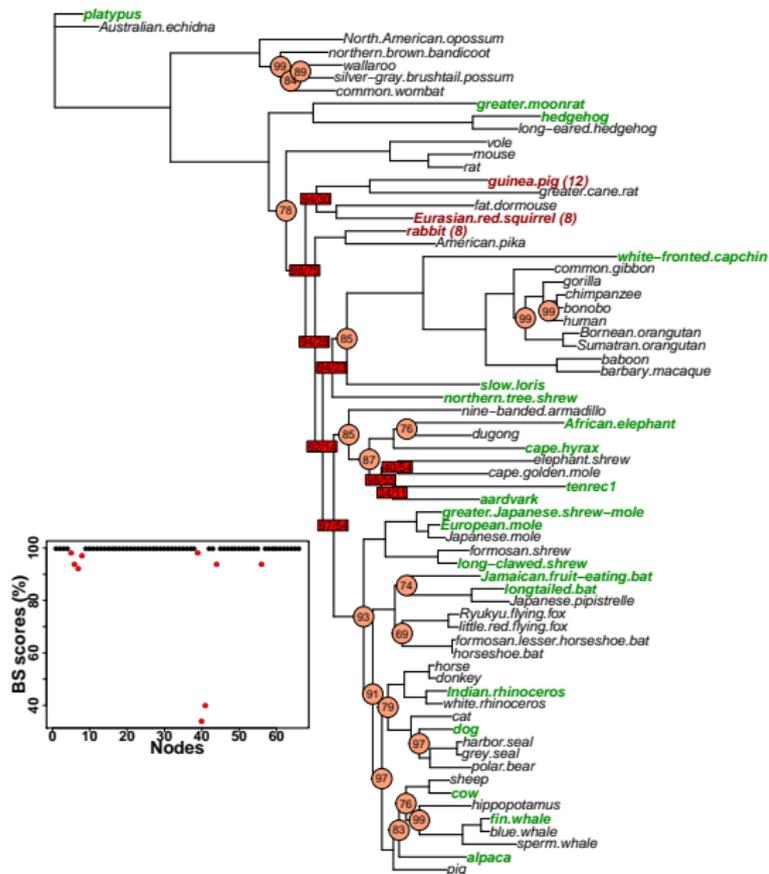


# Taxon Influence Index

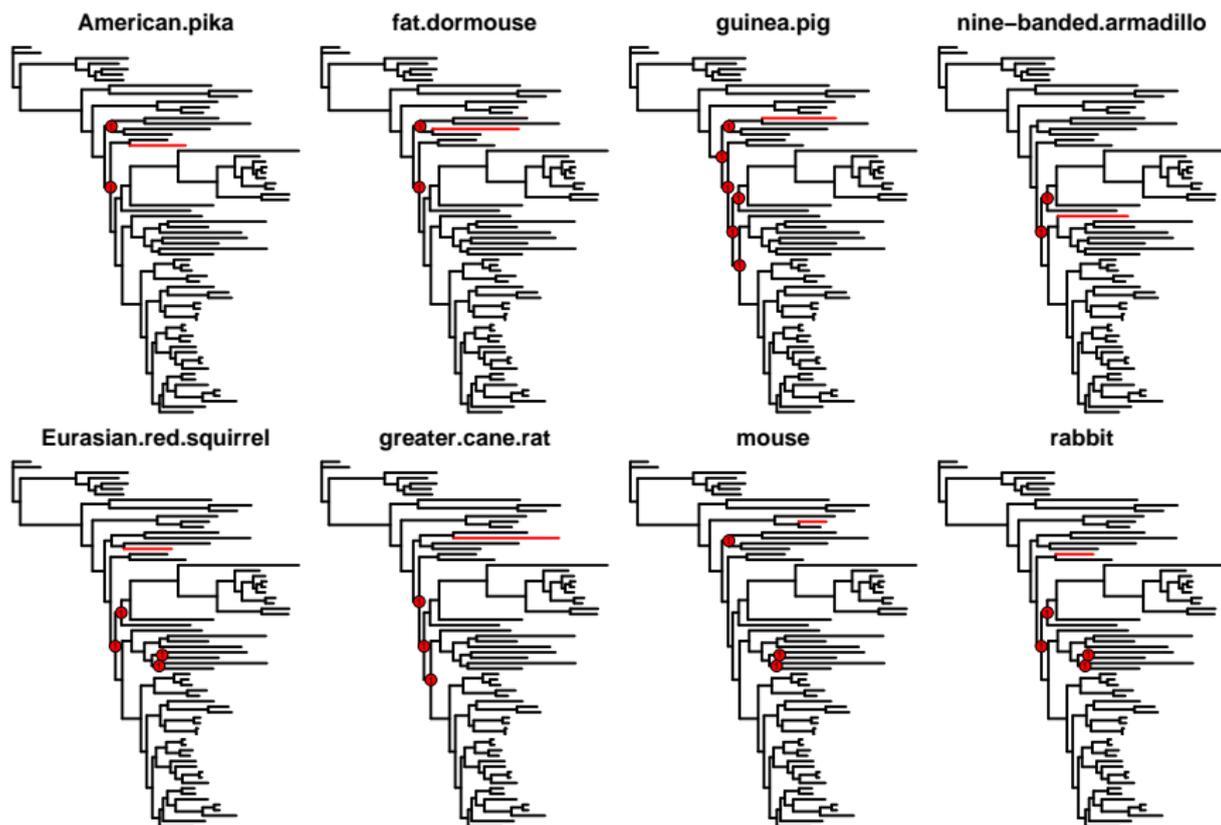


BSD: Branch-score distance / RF: Robinson-Foulds distance  
mean edge length: 0.05

# Complete Phylogeny

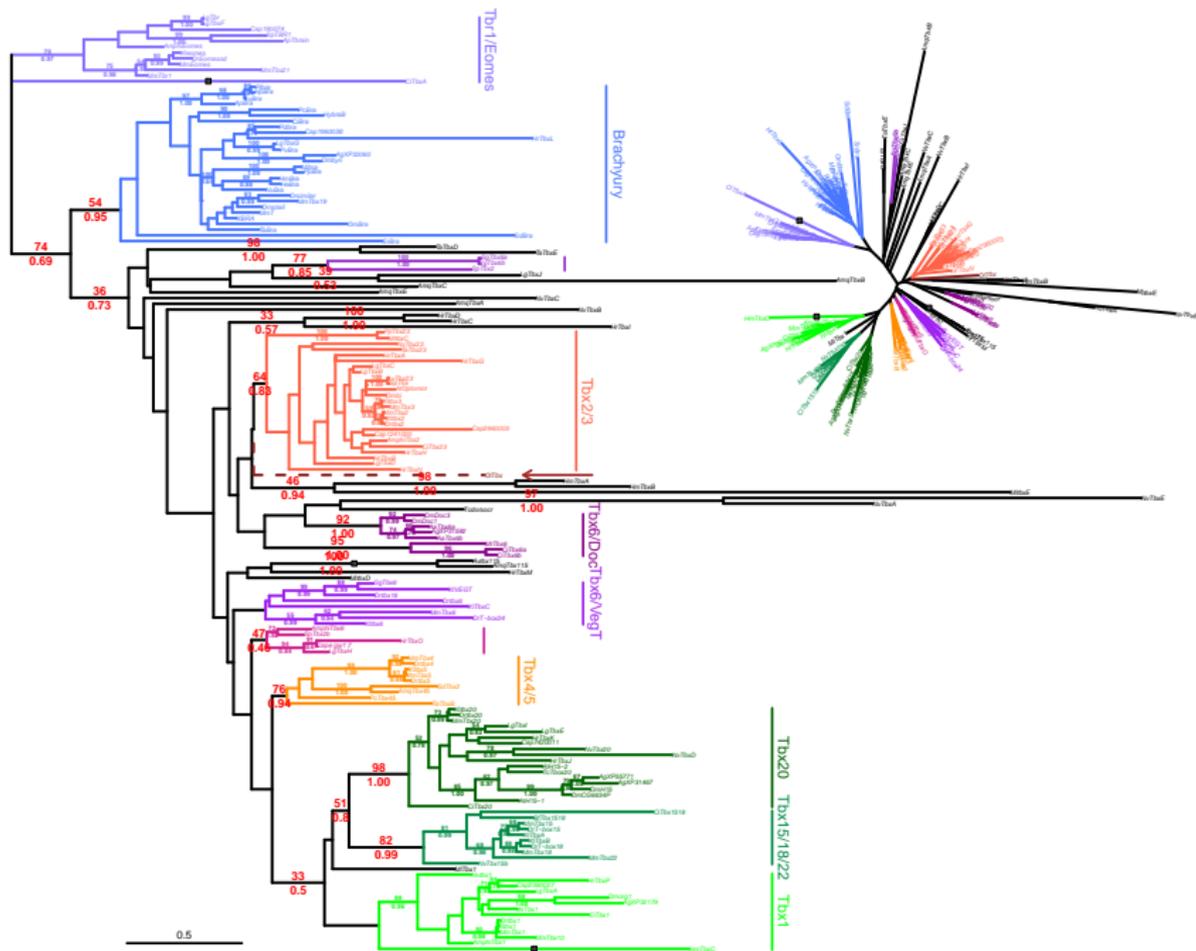


# Influential Species (mostly in Afrotheria)



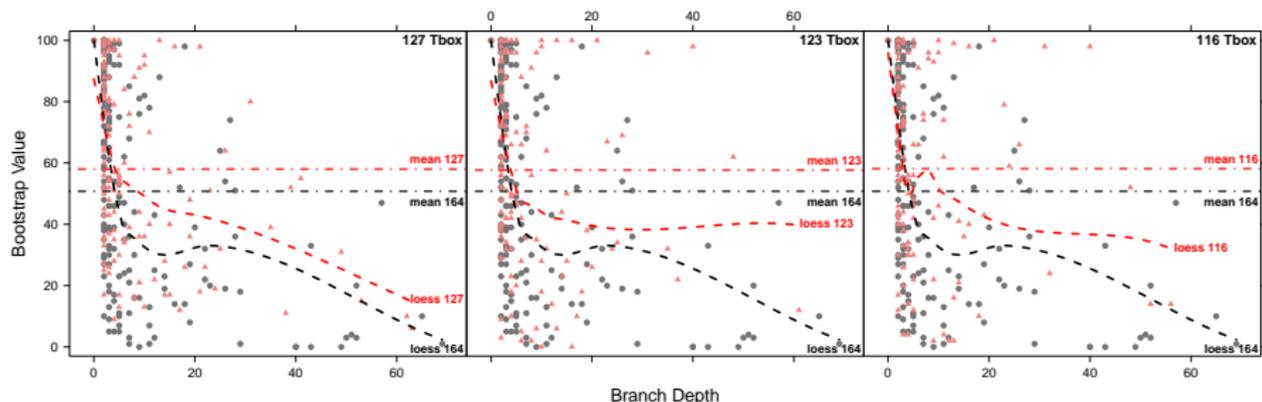
# Data: Bilaterian Transcription Factor T-Box

- T-box TF of 164 metazoans, involved in gastrulation;
- Amino Acids sequences;
- Sequences are 296 AA long;
- LG + I +  $\Gamma$ 4 model;
- Ancient family with 8 subfamilies
  - *Brachyury*
  - *Tbx1/10*
  - *Tbx15/18/22*
  - *Tbx20*
  - *Tbx2/3*
  - *Tbx4/5*
  - *Tbx6/VegT*
  - *Eomes/Tbr1/Tbx21*
- Interest lies in the position of TF **OITbx** present in a *Oscarella lobularis*.





# T-Box TF Phylogeny



Number of Tbx in tree	164	127	123	116
BP of OITbx in clade TBX2/3	19	32	34	59

No sponge (yet) identified as a member of the Tbx2/3 subfamily: **new evolutionary hypothesis.**

## Potential Sources of Instability and Indexes to Detect Them

- 1 Data sampling: Bootstrap;
- 2 Outliers: Influence function for sites;
- 3 Rogue species: Taxon Influence Index.

## Pros and Cons

- 1 Assess uncertainty coming from different sources;
- 2 Highlight potential outliers;
- 3 No rigorous statistical threshold for inclusion (p-values,...)