# Exploring Topological Incongruence for Detecting Contaminations in Phylogenomic Data Sets

## Frédéric Delsuc, Khalid Belkhir & Celine Scornavacca

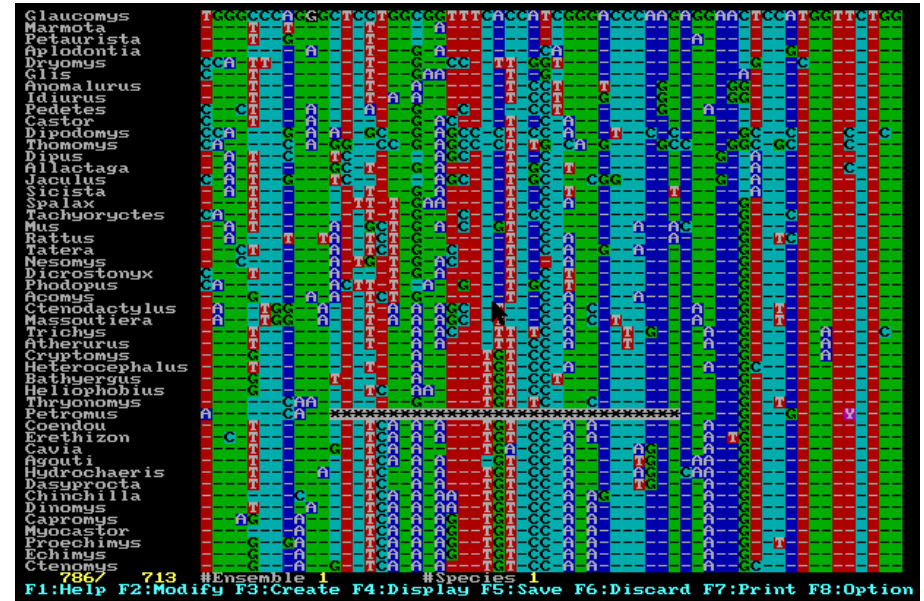**Institut des Sciences de l'Evolution - UMR 5554 - CNRS - IRD**
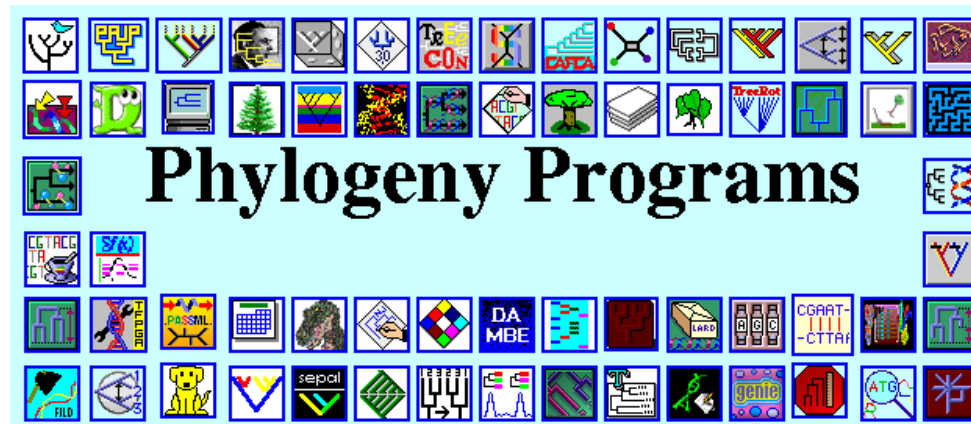**Université Montpellier 2 - France**

# Phylogenetic Reconstruction in Practice
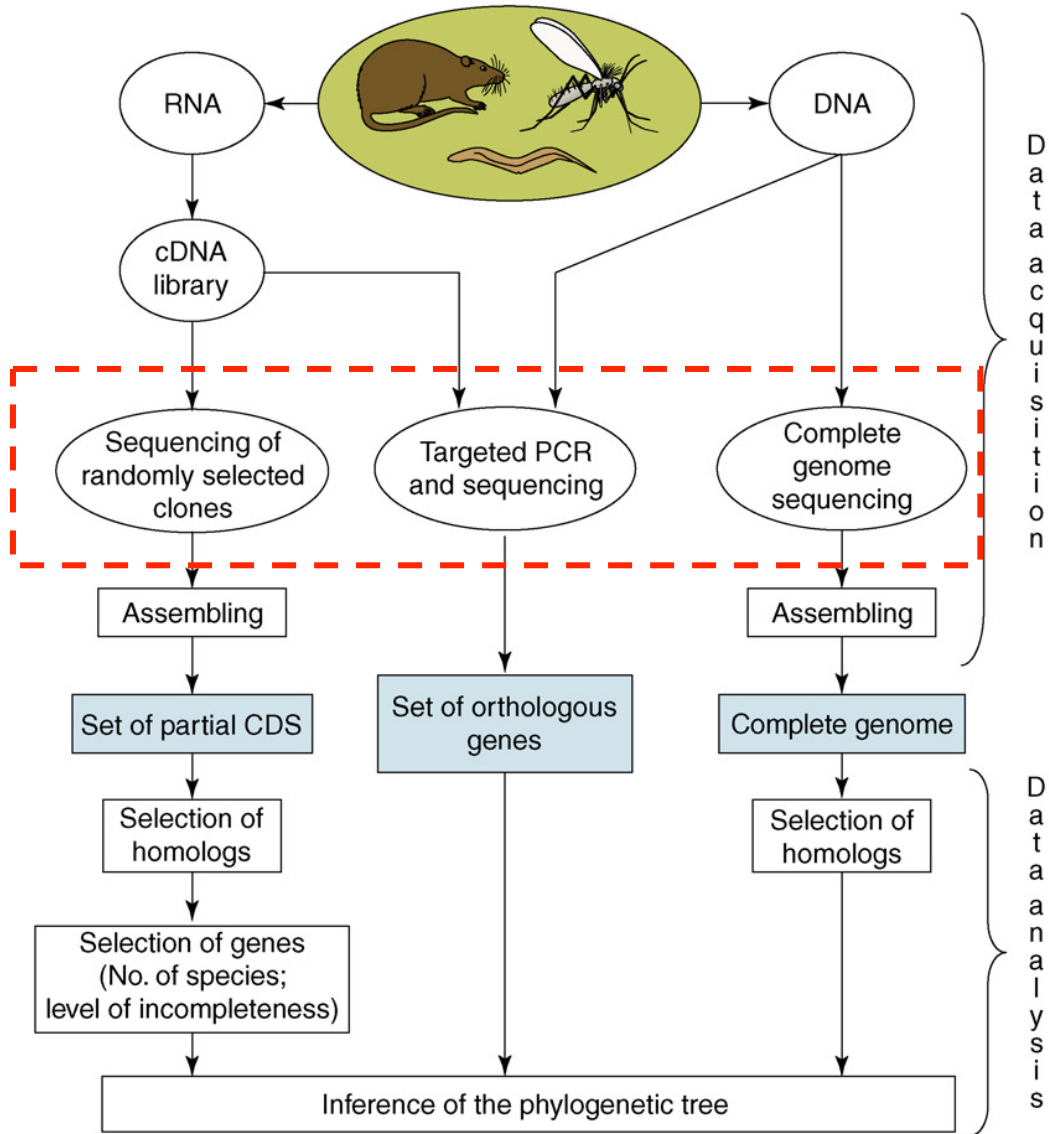


**1. Organisms**



**2. Homologous characters**

**3. Reconstruction methods**

# Data Sources in Phylogenomics



NGS

# Sequence Capture Methods for Phylogenomics



http://anchoredphylogeny.com/

512 nuclear loci
for vertebrates



UCEs identified in alignments of birds and lizard

http://ultraconserved.org/

up to 5,000 loci
in vertebrates

Lemmon *et al.* (2012) *Syst. Biol.*

Faircloth *et al.* (2012) *Syst. Biol.*

# The Future of Phylogenomics

# From Phylogenetics to Phylogenomics

# Towards a Full Resolution of the Tree of Life?

**Perspective**

## Resolving Difficult Phylogenetic Questions: Why More Sequences Are Not Enough

Hervé Philippe[1]*, Henner Brinkmann[1], Dennis V. Lavrov[2], D. Timothy J. Littlewood[3], Michael Manuel[4], Gert Wörheide[5,6], Denis Baurain[7]

**Sources of topological incongruence in phylogenomics:**

Biological reasons:
- Incomplete lineage sorting
- Horizontal gene transfer
- Hidden paralogy

Artificial reasons:
- Phylogenetic reconstruction artefacts (LBA, Compositional biases, …)
- **Sequence misidentifications / contaminations**

**=> Need for quality control methods and data exploration tools.**

# Phylogenomics and the Origin of Land Plants

**Multigene Phylogeny of the Green
Lineage Reveals the Origin
and Diversification of Land Plants**

*Coleochaete* **is the sister-group
of land plants**

77 nuclear genes
12,149 amino acid sites
77 plant taxa



Finet *et al.* (2010) *Current Biol.*

# Phylogenomics and the Origin of Land Plants

Coleochaetales
paraphyly



Finet *et al.* (2010) *Current Biol.*

# Cross-contaminations among New Transcriptomes

# Massive Contaminations in Finet et al. Data Set

We found a total of 101 contaminated sequences, including a rotifer instead of the charalean Nitella (rpl27), or a diatom instead of the chlorophyte Volvox (rpl11b); contaminations by parasites, symbionts or commensals are not rare in transcriptomic datasets ([5] and unpublished results) and should be systematically verified and discarded. More problematically, most (55 out of 101) correspond to cross-contaminations among the seven newly sequenced charophytes [4], i.e. sequences from distantly related charophytes are virtually identical at the nucleotide level. In particular, 29 sequences from the coleochaetalean Chaetosphaeridium are from the zygnematalean Penium.

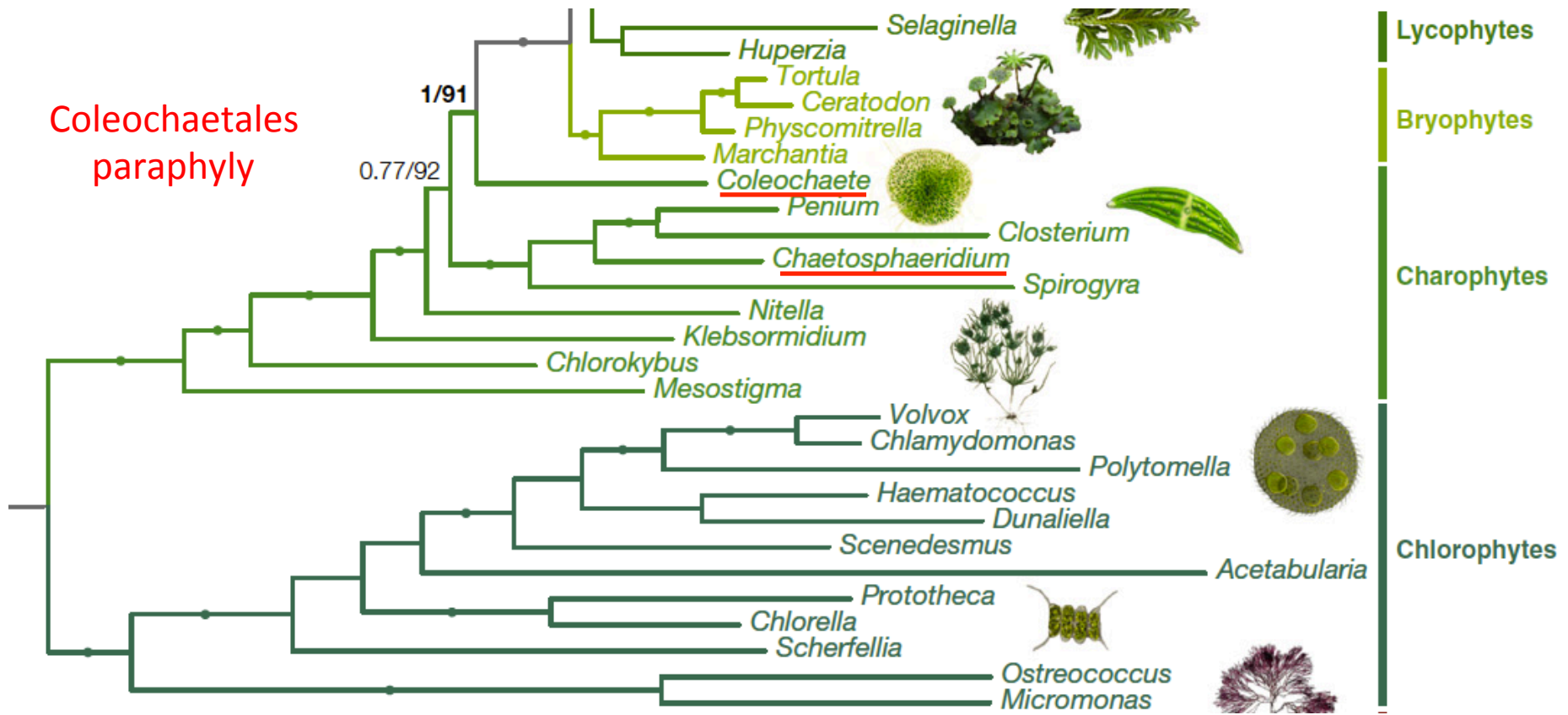| Contaminated | Contaminants | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| Studied organisms | Penium | Spirogy | Chaetos | Coeloch | Nitella | Klebsor | Chlorok | Non charophytes |
| Penium | | | | | | | 1(1) | 1(1) |
| Spirogyra | | | | | | | | 2(1) |
| Chaetosphaeridium | 29(29) | | | | | | 3(3) | 5(2) |
| Coleochaete | | 1(0) | | | | | | 1(1) |
| Nitella | 7(7) | | 5(2) | | | | | 5(3) |
| Klebsormidium | | | 6(5) | | | | 3(3) | 1(1) |
| Chlorokybus | | | | | | | | |
| Non Charophytes | | | | | | | | 31 (15) |

Laurin-Lemay, Brinkmann & Philippe (2012) *Current Biol.*

# Contaminations and the Origin of Land Plants

The congruence test revealed **74 contaminant sequences in the 77 ribosomal protein alignments** of Finet et al., and yielded to the **removal of 99 sequences** (because in 25 cases it was not possible to determine which is the correct sequence).



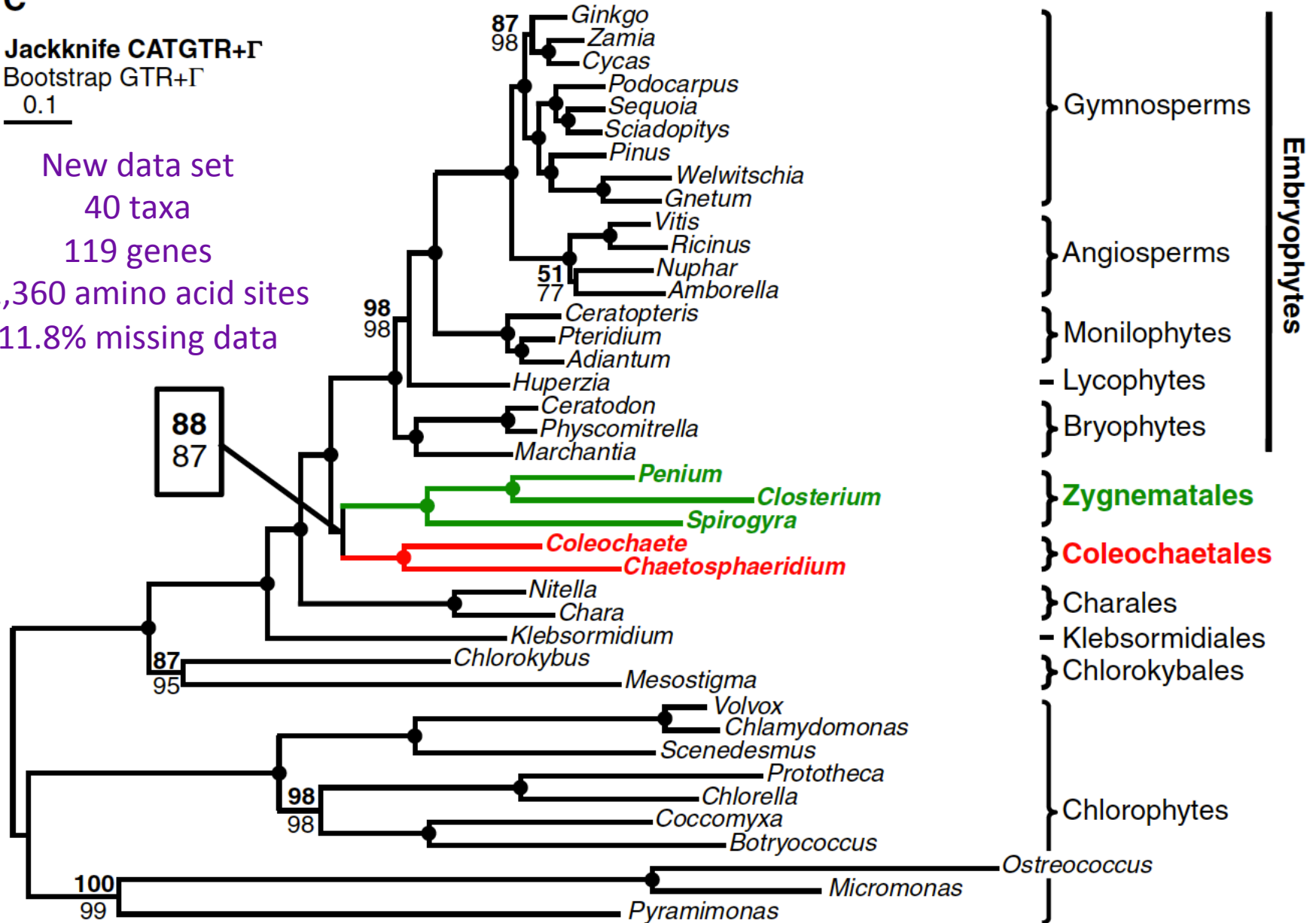**Original data set**                    **Decontaminated data set**

# Contamination-free Phylogenomics



Laurin-Lemay, Brinkmann & Philippe (2012) *Current Biol.*

# Methods for Detecting Outliers in Phylogenomic Data Sets

## Phylo-MCOA: A Fast and Efficient Method to Detect Outlier Genes and Species in Phylogenomics Using Multiple Co-inertia Analysis

Damien M. de Vienne,[*,1,2] Sébastien Ollier,[2] and Gabriela Aguileta[1,2]

=> Multiple co-inertia analysis (MCOA) extracting the similarities and discrepancies among genes in terms of pairwise distances.

## KDETREES: non-parametric estimation of phylogenetic tree distributions

Grady Weyenberg[1], Peter M. Huggins[2], Christopher L. Schardl[3], Daniel K. Howe[4] and Ruriko Yoshida[1,*]

=> Non-parametric method based on topological distances with the goal of identifying trees that are significantly different from the rest of distribution.

## TreSpEx—Detection of Misleading Signal in Phylogenetic Reconstructions Based on Tree Information

Torsten H. Struck

=> Combines different approaches utilizing tree-based information (nodal support or patristic distances) to identify misleading signals.

de Vienne *et al.* (2012) *Mol. Biol. Evol.*; Weyenberg et al. (2014) *Bioinformatics*; Struck (2014) *Evol. Bioinf.*

# A Biological Case Study: The Sister-group of Primates



Mining GenBank

17 genes

78 species

Lartillot & Delsuc (2012) *Evolution*.

# An Unexpected (yet Exciting) Result!



BP = 100 / PP =1.0

15,515
nucleotide sites
ML Tree
GTR+G8

*Trichechus*
*Loxodonta*
*Procavia*
*Orycteropus*
*Amblysomus*
*Echinops*
*Elephantulus*
*Macroscelides*
*Choloepus*
*Myrmecophaga*
*Tamandua*
*Dasypus*
*Euphractus*
*Chaetophractus*
*Solenodon*
*Sorex*
*Erinaceus*
*Galemys*
*Talpa*
*Megaderma*
*Rousettus*
*Pteropus*
*Nycteris*
*Artibeus*
*Tadarida*
*Myotis*
*Antrozous*
*Vicugna*
*Lama*
*Sus*
*Tragelaphus*
*Bos*
*Hippopotamus*
*Tursiops*
*Megaptera*
*Equus*
*Ceratotherium*
*Tapirus*
*Manis*
*Canis*
*Ailuropoda*
*Panthera*
*Felis*
*Otolemur*
*Microcebus*
*Lemur*
*Tarsius*
*Callithrix*
*Macaca*
*Pongo*
*Gorilla*
*Pan*
*Homo*

**Primates**

*Galeopterus*
*Cynocephalus*   **Dermoptera**    **SUNDATHERIA**
*Ptilocercus*   **Scandentia**
*Tupaia*
*Urogale*
*Ochotona*
*Sylvilagus*
*Oryctolagus*
*Muscardinus*
*Tamias*
*Spermophilus*
*Hystrix*
*Erethizon*
*Hydrochoerus*
*Cavia*
*Dipodomys*
*Castor*
*Pedetes*   **Glires**
*Mus*
*Rattus*

0.06

# Support for Primatomorpha



Primatomorpha
ML BS: 90%
BAY PP: 1.00
Indels:
2 aa del., SPBC25
2 aa del., SMPD3
4 aa del., MTUS1
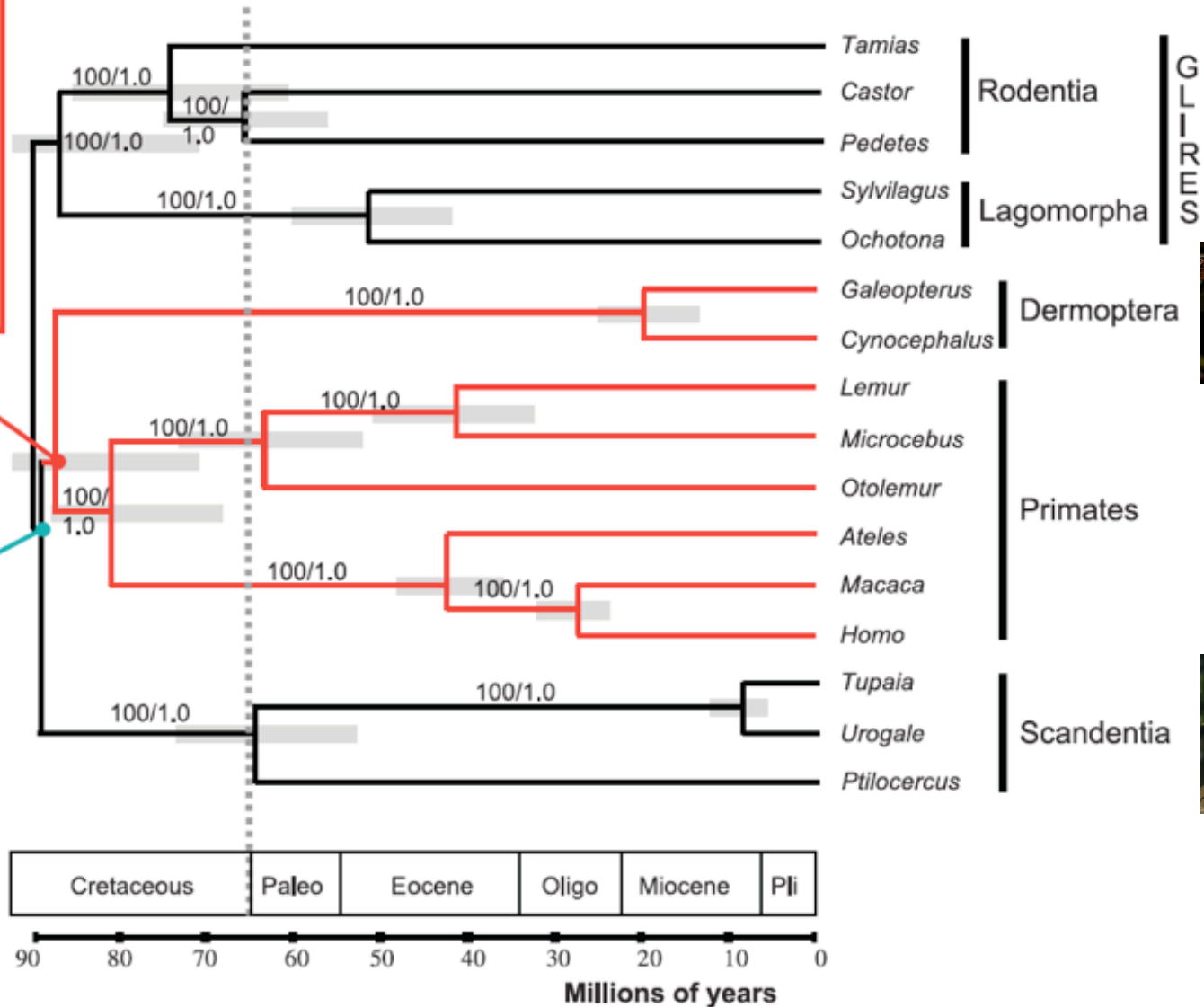3 aa del., SH3RF2
4 aa del., NCOA4
3 aa del., TEX2
1 aa del., SSH2

Euarchonta
ML BS: 92%
BAY PP: 1.00
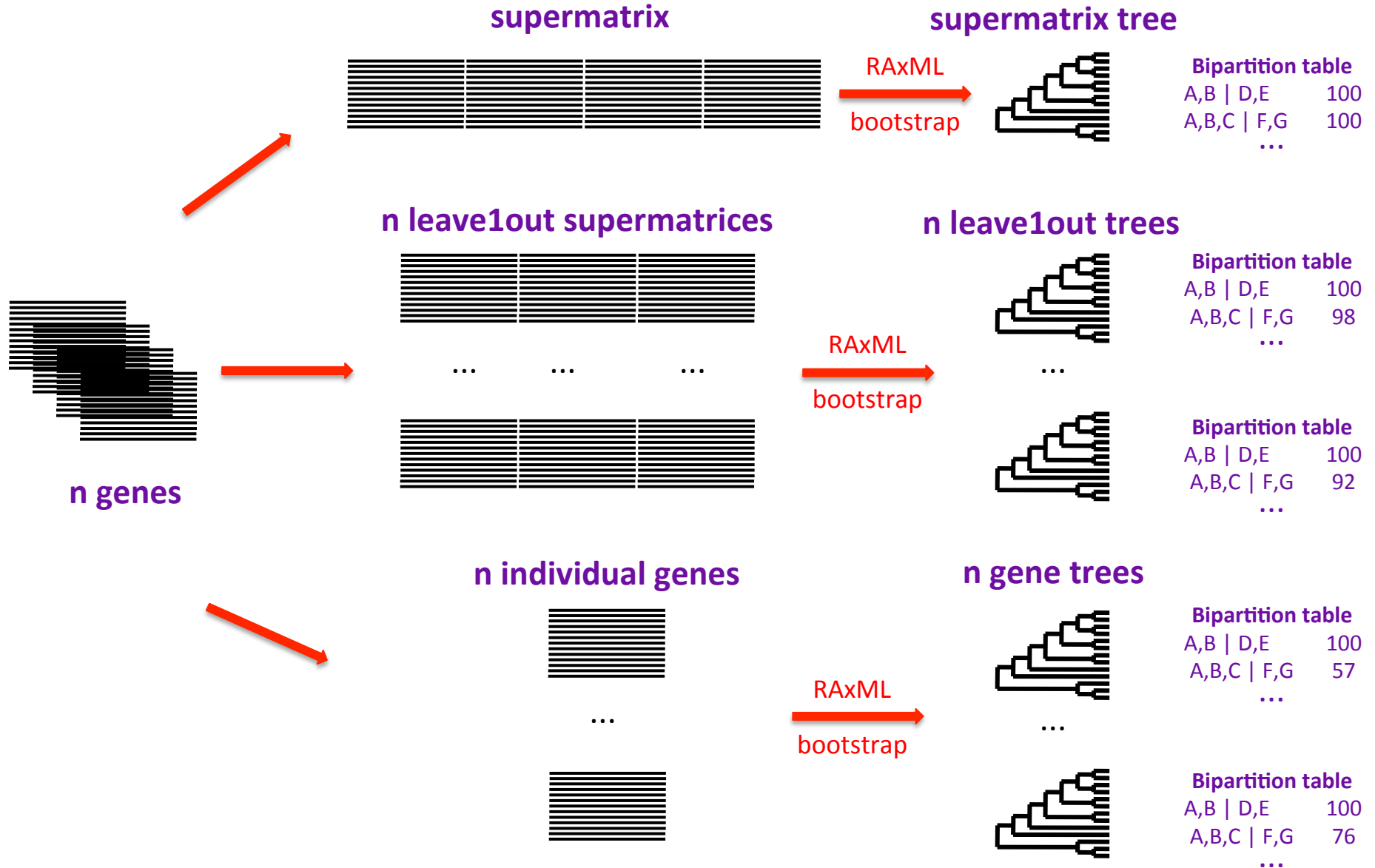Indels:
4 aa del, N4BP2
3 aa ins., ZNF12
1 aa del., CDCA5

Janecka *et al.* (2007) *Science*

# A Simple Pipeline for Exploring Topological Incongruence



**supermatrix**

**supermatrix tree**

RAxML
bootstrap

**Bipartition table**
A,B | D,E          100
A,B,C | F,G        100
...

**n leave1out supermatrices**

**n leave1out trees**

RAxML
bootstrap

**Bipartition table**
A,B | D,E          100
A,B,C | F,G         98
...

**Bipartition table**
A,B | D,E          100
A,B,C | F,G         92
...

**n genes**

**n individual genes**

**n gene trees**

RAxML
bootstrap

**Bipartition table**
A,B | D,E          100
A,B,C | F,G         57
...

**Bipartition table**
A,B | D,E          100
A,B,C | F,G         76
...

# Montpellier Bioinformatics Biodiversity Platform (MBB)

http://mbb.univ-montp2.fr/MBB/

# Exploratory Plots for Detecting Incongruent Bipartitions



**MBB** — Login

NGS   Phylogenomics   Population Genetics   Population Dynamics   Ecological Modelling

All Online Tools   All Downloads   All Data   All Docs & Events   Search

**Online Tools**
- Blast my DB
- Ima
- PhyML
- MAFFT
- Structure
- more ...

**Dowloads**
- Genetix
- Migraine
- CoMap
- Genepop
- HGT simul
- more ...

**Data**
- Polymorphix
- OrthoMaM
- more ...

**Misc**
- Platform Load
- Pubmed RSS feed
- Other services
- FAQ
- Contact

*SUNDAJOY*

**Results:**

trees/genes_vs_SM.txt (37.70 Ko)

trees/SM_vs_L1O.txt (9.42 Ko)

trees/SM_vs_genes.txt (8.26 Ko)

trees/RAXML_bestTrees.nex (80.91

trees/SM_vs_genes.html (2.33 Ko)

trees/SM_vs_L1O.html (2.33 Ko)

trees/genes_vs_SM.html (2.33 Ko)

SUNDAJOY.out (2.65 Ko)

standard error file

*Unix exact command:*

SUNDAJOY Sundacont.zip FILENAMES PWD 100 /share/apps/bin/RAxML-7.2.8-ALPHA/ Fasta SGE 90

*Your input data:*

Sundacont.zip

*Pise CGI generator*

**1. Support for Supermatrix bipartitions in Leave1out trees**

**2. Support for Supermatrix bipartitions in Gene trees**

**3. Support for Gene tree bipartitions not in the Supermatrix**

# Exploratory Plots for Detecting Incongruent Bipartitions

# Contaminations / Misindentifications in GenBank



Nycteris (Chiroptera) / Rattus (Rodentia)

**TYR (Tyrosinase)**



BP = 100

Problem:
*Nycteris grandis* (AY834610.1)
is a murid rodent

# Contaminations / Misindentifications in GenBank


Sorex (Eulipotyphla) / Sus (Cetartiodactyla)

**PNOC (Prepronociceptin)**



BP = 100

Problem:
*Sorex* araneus (AY011813.1)
is a Suidae

# Contaminations / Misindentifications in GenBank

**APOB (Apolipoprotein B)**



BP = 100

*Tupaia* is a Dermopteran (97% similarity)

>gi|256549376|gb|FJ648363.1| *Tupaia glis* apolipoprotein B-like (APOB) gene
Submitted 29 October 2008
Ali F, Pons J, Shekelle M, Goodman M and Meier R
A sparse supermatrix recovers a well-supported primate phylogeny with dates (Unpublished)

# Effects of Contaminations on Phylogenomic Inference



**SUNDATHERIA**

**PRIMATOMORPHA**

**Original data set**

**Decontaminated data set**

0.06

0.06

# Effects of Contaminations on Phylogenomic Inference



**Original data set**

**Decontaminated data set**

SUNDATHERIA

PRIMATOMORPHA

0.06

# Conclusions

=> **Contaminations / misidentifications are frequent** in phylogenomic data sets and public databases

=> Simple **data exploration tools** based on bipartition support allow detecting incongruent signals due to misidentified / cotanminated sequences

=> Even a few number of contaminations / misidentifications **can strongly impact phylogenetic inference** when phylogenetic signal is scarce

# Thanks for your attention!

# Applications of Sequence Capture Methods



Lemmon & Lemmon (2013) *Annu. Rev. Ecol. Evol. Syst.*

SUNDACLEAN

*Choloepus*
*Tamandua*
*Myrmecophaga*
*Dasypus*
*Chaetophractus*
*Euphractus*
*Trichechus*
*Loxodonta*
*Procavia*
*Orycteropus*
*Echinops*
*Amblysomus*
*Macroscelides*
*Elephantulus*
*Solenodon*
*Talpa*
*Galemys*
*Erinaceus*
*Sorex*
*Megaderma*
*Pteropus*
*Rousettus*
*Nycteris*
*Artibeus*
*Tadarida*
*Antrozous*
*Myotis*
*Vicugna*
*Lama*
*Sus*
*Tragelaphus*
*Bos*
*Hippopotamus*
*Tursiops*
*Megaptera*
*Equus*
*Ceratotherium*
*Tapirus*
*Manis*
*Panthera*
*Felis*
*Ailuropoda*
*Canis*
*Galeopterus*
*Cynocephalus*
*Otolemur*
*Lemur*
*Microcebus*
*Tarsius*
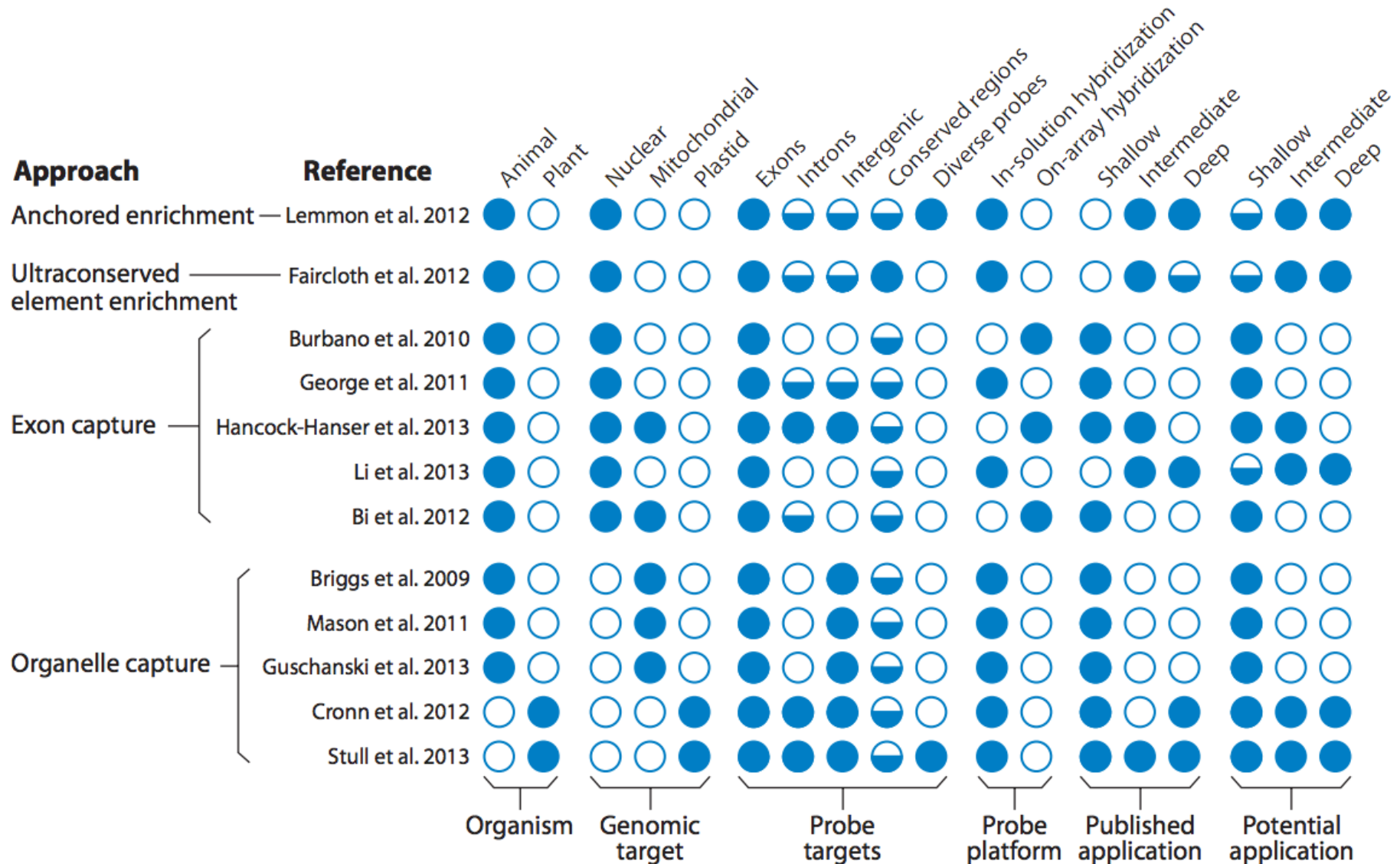*Callithrix*
*Macaca*
*Pongo*
*Gorilla*
*Homo*
*Pan*
*Ptilocercus*
*Urogale*
*Tupaia*
*Ochotona*
*Sylvilagus*
*Oryctolagus*
*Muscardinus*
*Spermophilus*
*Tamias*
*Hystrix*
*Erethizon*
*Hydrochoerus*
*Cavia*
*Castor*
*Dipodomys*
*Pedetes*
*Mus*
*Rattus*

**Dermoptera**

**PRIMATOMORPHA**

**Primates**

**Scandentia**

**Glires**

0.06

# A Jackknife Approach Based on the Supermatrix



Legend:
- (Lion - Leopard - Jaguar) (Tiger - Snow Leopard - Clouded Leopard)
- Lion - Leopard
- Lion - Jaguar
- Tiger - Snow Leopard
- Snow Leopard - Clouded Leopard
- Leopard - Tiger

Y-axis: Nonparametric Bootstrap Support

X-axis: ALL, FGB, IRBP, TTR, APP, CALB, CHRNA, CLU, CMA, DGKG2, FES, GATA, GHR, GNAZ, GNB, HK1, NCL, PNOC, RAG2, RASA, SILV, TCP

# Methods of Phylogenomic Inference



Delsuc, Brinkmann & Philippe (2005). *Nat. Rev. Genet.* 6: 361-375.

# Phylogenomics Increases the Resolving Power



Delsuc, Brinkmann & Philippe (2005) *Nat. Rev. Genet.*