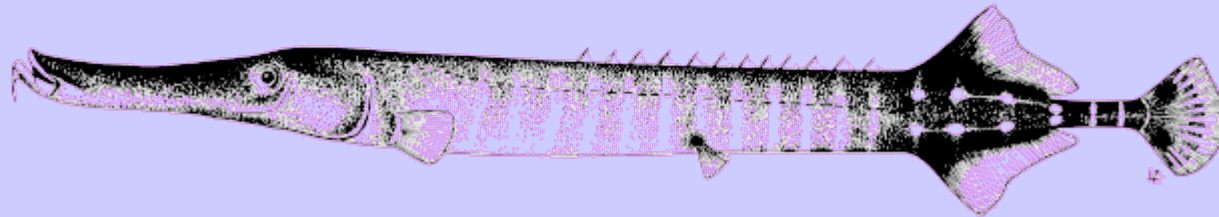# Towards a Reliability Index for Clades: an Application on Acantomorph Teleosts

Blaise Li and Guillaume Lecointre

UMR 7138
Département Systématique et Évolution
Muséum National d'Histoire Naturelle - Paris

ASIH annual meeting - 11 July, 2005
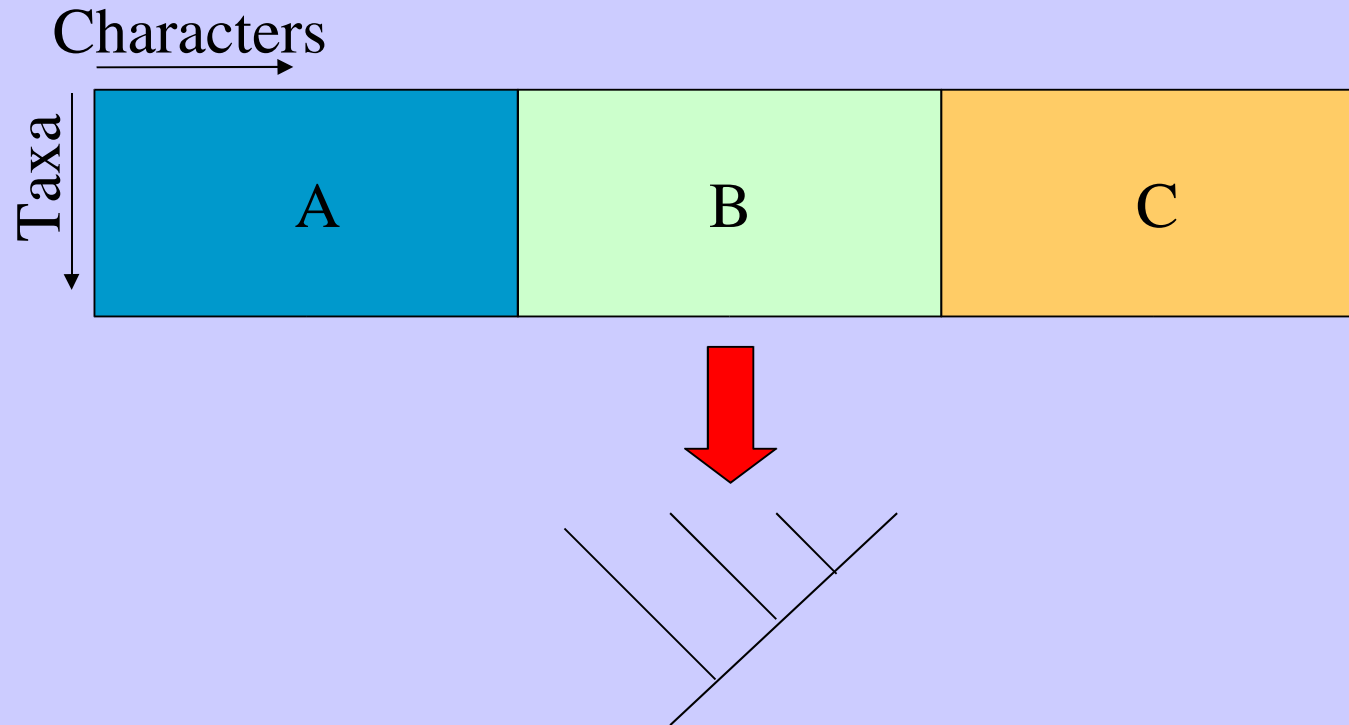
# How to measure a result's quality ?

- Does it resist data perturbation (robustness) ?

- What is the credit given to a statement about relationships among species (**reliability**) ?

# How to measure reliability ?
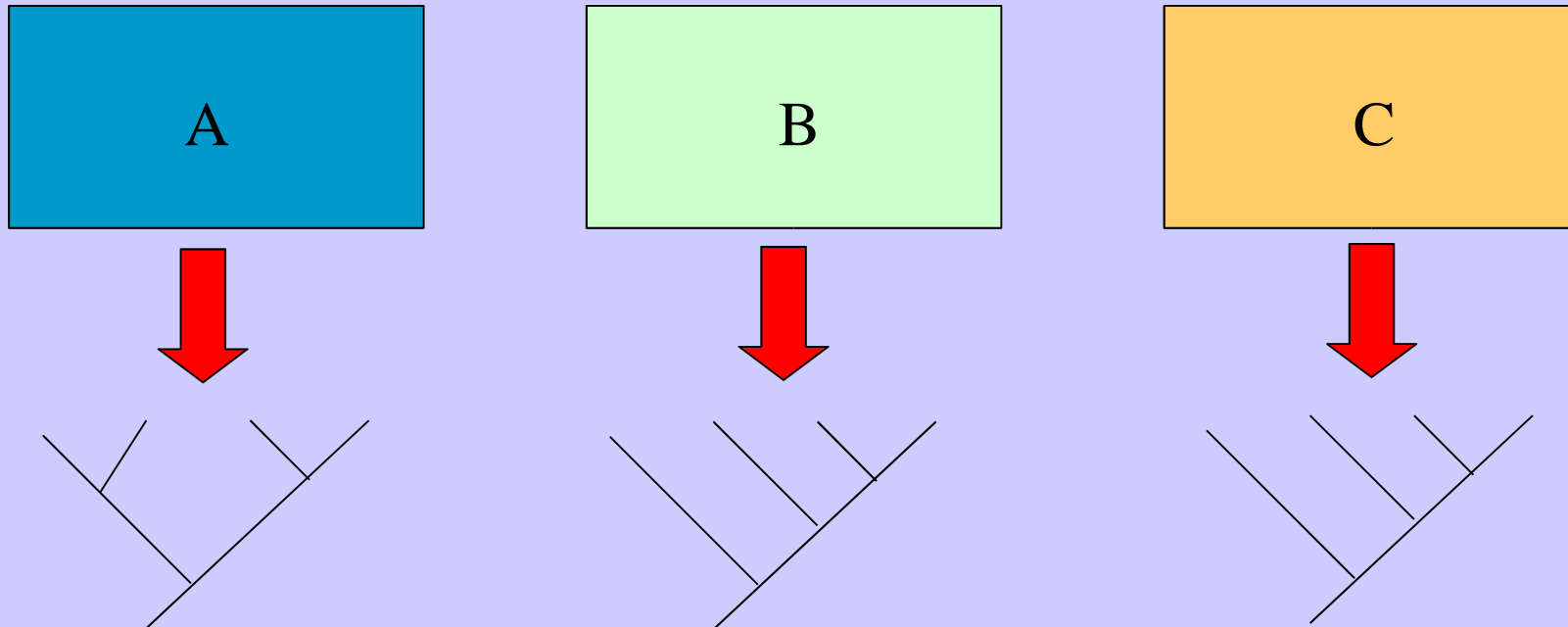
Use multiple data sets

Combine the data into a single matrix or keep it separate ?

# Combined analysis (« Total Evidence »)



Drawback:    a marker-specific bias can influence the inference from the whole data during the optimization process

# Separate analysis



Biological background knowledge is needed to justify delineation and independence of the data sets

Drawbacks:  higher stochastic effects
full expression of marker-specific biases

# How to measure reliability ?

Consider corroboration between **independent parts** of the data (with **partial data combination**)
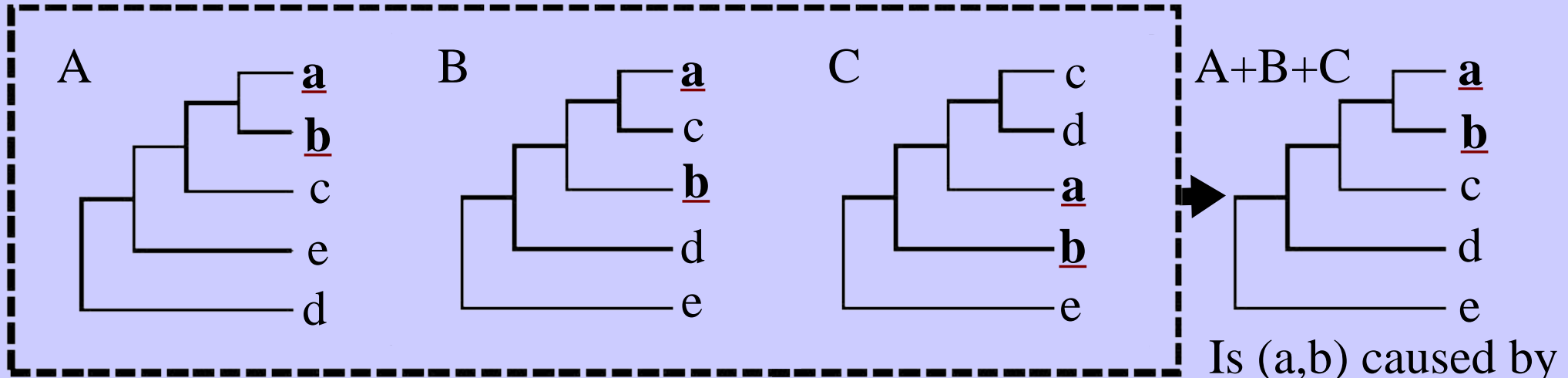
# About independent data

Choose basic markers you postulate to evolve independently

Thus, if a marker is subject to some bias, you would expect that bias not to exist for the other markers also

It increases the chance that **repeated results** are caused by a shared feature of the data sets: the history of the taxa represented by the markers

# Partial combinations...
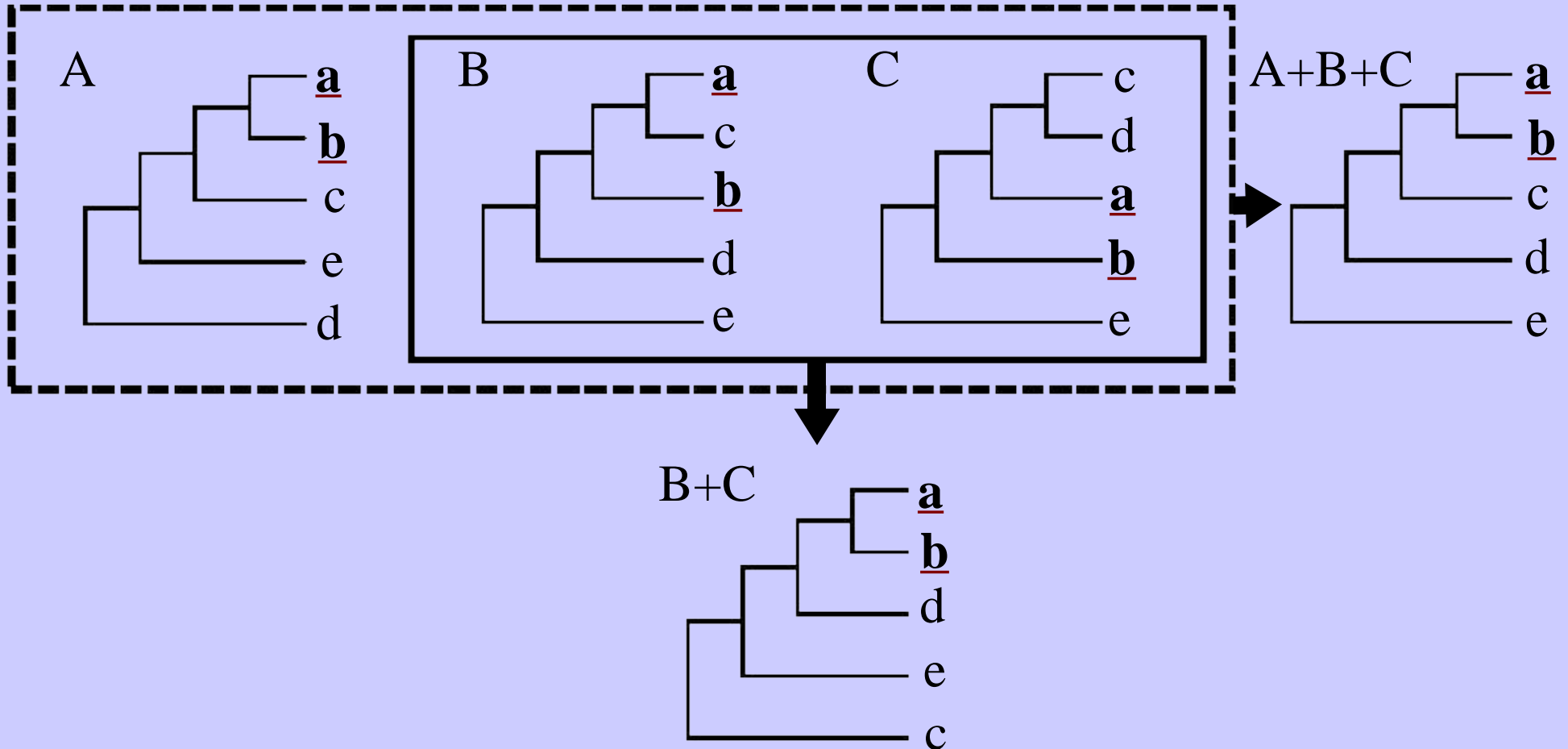


A bears signal for (a,b). Is it true historical signal ?

Is (a,b) caused by A only ?

# ... to reduce stochastic effects



B+C corroborates (a,b). Partial combination overcame some stochastic effects that prevented (a,b) from appearing.

# A repeatability index for clades: a way to formalize reliability

The more a clade is found by the analysis of independent data, the more reliable it is

- Separate the data into non-overlapping **parts** and analyze each part with the same method

**data sets: A, B, …**

A    B    1

E    C

3    D    2

One of the possible partitioning schemes with 3 parts (1, 2, 3)

# A repeatability index for clades: a way to formalize reliability

The more a clade is found by the analysis of independent data, the more reliable it is

- Separate the data into non-overlapping **parts** and analyze each part with the same method

- Count the occurrences of the clades among non-overlapping parts

- Repeat the process with each possible **partitioning scheme**

- For each clade, retain the highest number of occurrences over all the partitioning schemes

**A partitioning scheme with 3 parts**

data sets: A, B, ...

| clade | occurrences |
|-------|-------------|
| $\alpha$: | **3** occurrences |
| $\beta$: | **3** occurrences |
| $\gamma$: | **2** occurrences |
| $\delta$: | **2** occurrences |
| $\epsilon$: | 1 occurrence |
| $\zeta$: | 1 occurrence |
| ... | |

**Another partitioning scheme with 3 parts**

data sets: A, B, ...

| clade | occurrences |
|-------|-------------|
| $\alpha$: | 2 occurrences |
| $\beta$: | **3** occurrences |
| $\gamma$: | 1 occurrence |
| $\delta$: | **2** occurrences |
| $\epsilon$: | **3** occurrences |
| $\zeta$: | **2** occurrences |
| ... | |

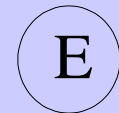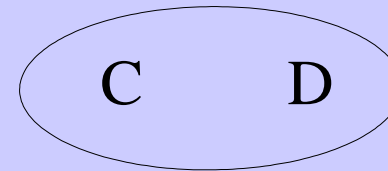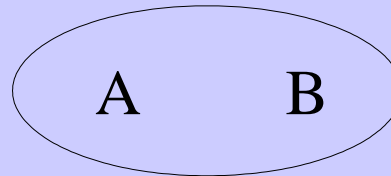# Taking into account probabilities for clades to appear at random

- Clades appear also from data containing no signal

- Therefore one must measure the probability of a clade to appear purely at random for a given program:

  - give random data to the program

  - analyze that data

  - repeat the process and compute the frequency of each clade

**data sets**

|  | part 1 | | part 2 | | part 3 |
|---|---|---|---|---|---|
|  | ( A    B ) | | ( C    D ) | | ( E ) |
| $\alpha$ (3 occurrences): | $1-P_1(\alpha)$ | $+$ | $1-P_2(\alpha)$ | $+$ | $1-P_3(\alpha)$ |
| $\beta$ (3 occurrences): | $1-P_1(\beta)$ | $+$ | $1-P_2(\beta)$ | $+$ | $1-P_3(\beta)$ |
| $\gamma$ (2 occurrences): | $1-P_1(\gamma)$ | $+$ | $0$ | $+$ | $1-P_3(\gamma)$ |
| $\delta$ (2 occurrences): | $0$ | $+$ | $1-P_2(\delta)$ | $+$ | $1-P_3(\delta)$ |
| $\epsilon$ (1 occurrence): | $1-P_1(\epsilon)$ | $+$ | $0$ | $+$ | $0$ |
| $\zeta$ (1 occurrence): | $0$ | $+$ | $1-P_2(\zeta)$ | $+$ | $0$ |
| ... | ... | | ... | | ... |

# Taking into account contradiction among clades

- First order reliability of clade $\alpha$:

$$R_1(\alpha)=\text{Max}_{\text{partitioning schemes}}(\sum_{\text{parts}}(\text{Occurrences-P}))$$

- Second order reliability of clade $\alpha$:

$$R_2(\alpha)=R_1(\alpha)-R_1(\beta_1)$$

$\beta_1$ being the highest $R_1$ $\alpha$ contradictor

- Etc...

$$R_n(\alpha)=R_1(\alpha)-R_1(\beta_{n-1})$$

# An application on Acanthomorpha

- 5 basic data sets:

  – mitochondrial markers: partial 12S+16S (828 bp)

  – nuclear markers: partial 28S (801 bp), partial Rhodopsin (759 bp), partial MLL (552 bp), partial IRBP (713 bp)

- 73 taxa shared by all data sets

- 31 ways to combine the 5 independent data sets (allowing to compose 51 partitioning schemes), each analyzed with Paup 4 under Maximum Parsimony

Total evidence majority rule consensus tree

Majority-rule bootstrap tree for the total data combination

Tree constructed to include the highest reliability inter-compatible clades

# A few reliable clades



5/4.9999
A
   5/4.9980
— Zeus
— Zenopsis
— Trachyrincus

3/2.0078
K
   5/4.9980
— Perca
— Gymnocephalus
   4/3.0062
— Notothenia
— Bovichtus

3/2.0078
F
   5/4.9980
— Monopterus
— Mastacembelus
   5/4.9980
— Ctenopoma
— Channa

3/2.0077
G
— Uranoscopus
   5/4.0059
— Cheimarrichthys
— Ammodytes

# Reliability is not robustness

high reliability, neither present in the total evidence tree, nor in the bootstrap consensus



Total evidence



Best repeatability scores



$$\frac{\textbf{max occurrences}/\text{repeatability score}}{\textbf{bootstrap value}}$$

Bootstrap

# Reliability is not robustness



Total evidence

Best repeatability scores

max occurrences/repeatability score
bootstrap value

high robustness,
low reliability

Bootstrap

# A short conclusion because it's time to finish

- Some non-robust clades have been identified as reliable by our method (Psenopsis, Pampus)

- Some clades with high boostrap support are not considered reliable by our method (Dactylopterus, Aulostomus)

- Works with completely shared taxonomic samplings

# Acknowledgements

- G. Lecointre, W.-J. Chen and A. Dettaï

- Service de Systématique Moléculaire (MNHN)

- École Normale Supérieure and French Ministry of Research

- Free software programers

Regalecus
5/4.9980 ┌ Zeus
5/4.9999 └ Zenopsis
3/1.0156 Trachyrincus
Polymixia
4/3.0062 Beryx
Barbourisia
Myripristis
3/2.0066 ┌ Psenopsis
Pampus
5/4.9999 3/1.0145 ┌ Scomber
Kali
5/4.9999 Dactylopterus
2/0.0156 5/4.9999 3/1.0152 ┌ Macroramphosus
Aulostomus
Liza
2/0.0156 3/1.0152 ┌ Belone
Bedotia
1/0.0078 4/3.0063 ┌ Parablennius
1/-0.984 3/2.0078 Forsterygion
5/4.9980 ┌ Lepadogaster
Apletodon
5/4.9980 ┌ Monopterus
Mastacembelus
3/2.0078 5/4.9980 ┌ Ctenopoma
Channa
5/4.9999 Poecilia
1/0.0078 3/1.0152 ┌ Callionymus
2/0.0234 Mullus
Lagocephalus
Arnoglossus
Lateolabrax
2/0.5155 5/4.9980 ┌ Scarus
Labrus
Pholis
5/4.9999 5/4.9980 ┌ Taurulus
1/0.0078 3/1.0156 Cyclopterus
Chelidonichthys
Trachinus
2/0.0156 2/0.0148 ┌ Serranus
1/-0.984 2/0.0155 Scorpaena
1/0.0078 2/0.0156 Holanthias
5/4.9980 ┌ Rypticus
1/-0.492 2/0.0156 2/1.0077 Pogonoperca
Epinephelus

| Node label | Taxon |
|---|---|
| | Mullus |
| 2/0.0234 | Lagocephalus |
| | Arnoglossus |
| 2/0.5155 | Lateolabrax |
| 5/4.9980 | Scarus |
| | Labrus |
| 5/4.9999 | Pholis |
| 1/0.0078 | Taurulus |
| 3/1.0156 | Cyclopterus |
| 5/4.9980 | Chelidonichthys |
| | Trachinus |
| 2/0.0148 | Serranus |
| | Scorpaena |
| 2/0.0155 | Holanthias |
| 2/0.0156 | Rypticus |
| 2/1.0077 | Pogonoperca |
| | Epinephelus |
| 5/4.9980 | Perca |
| 3/2.0078 | Gymnocephalus |
| 4/3.0062 | Notothenia |
| | Bovichtus |
| 3/2.0077 | Uranoscopus |
| 5/4.0059 | Cheimarrichthys |
| | Ammodytes |
| 3/1.0144 | Ctenochaetus |
| | Capros |
| 3/1.0152 | Triacanthodes |
| | Ostracion |
| 5/4.9999 | Mola |
| | Chaetodon |
| 3/2.0077 | Ceratias |
| 4/3.0062 | Antennarius |
| | Holacanthus |
| | Pomadasys |
| | Mene |
| 3/1.5144 | Microchirus |
| | Citharus |
| 1/0.0077 | Psettodes |
| | Pentanemus |
| | Sphyraena |
| 3/2.0077 | Trachinotus |
| 3/2.0066 | Chloroscombrus |
| | Echeneis |
| | Halobatrachus |
| | Bathypterois |

# Reliability is not robustness (2)

**max occurrences**/repeatability score

**bootstrap value**



| | |
|---|---|
| 3/2.0077 | Chaetodon |
| 4/3.0062 | Ceratias |
| 52.8 | |
| 95.0 | Antennarius |

Not very robust
but rather reliable

Bootstrap

# Limits of the method



2/0.0234

1/0.0078 — Poecilia
3/1.0152 — Callionymus
— Mullus
— Lagocephalus
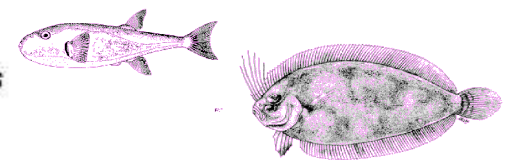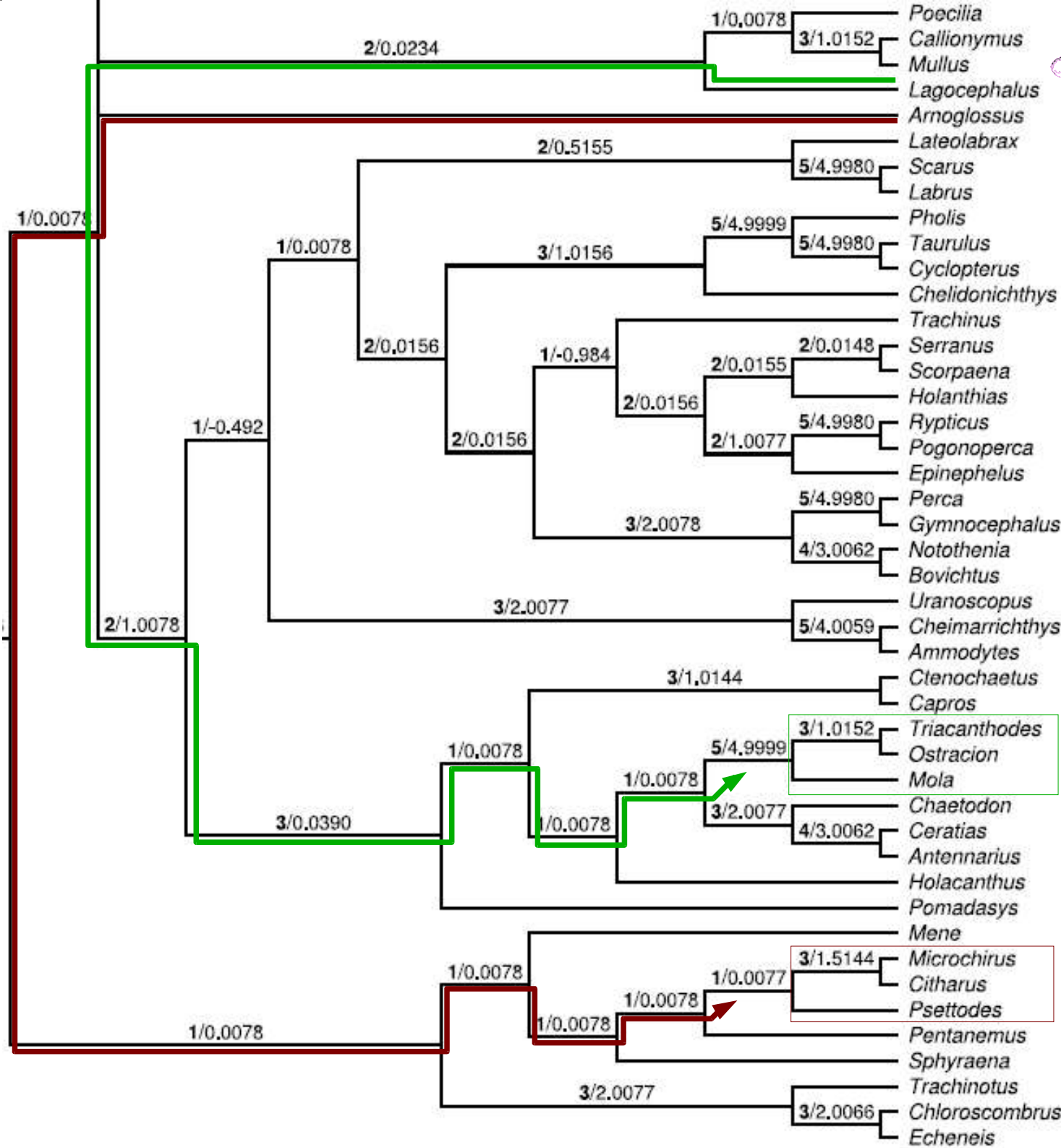— Arnoglossus

General long-branch
attraction is still a problem

Best repeatability scores

Long branches to swim back, following the low-reliability track !

Poecilia
Callionymus
Mullus
Lagocephalus
Arnoglossus
Lateolabrax
Scarus
Labrus
Pholis
Taurulus
Cyclopterus
Chelidonichthys
Trachinus
Serranus
Scorpaena
Holanthias
Rypticus
Pogonoperca
Epinephelus
Perca
Gymnocephalus
Notothenia
Bovichtus
Uranoscopus
Cheimarrichthys
Ammodytes
Ctenochaetus
Capros
Triacanthodes
Ostracion
Mola
Chaetodon
Ceratias
Antennarius
Holacanthus
Pomadasys
Mene
Microchirus
Citharus
Psettodes
Pentanemus
Sphyraena
Trachinotus
Chloroscombrus
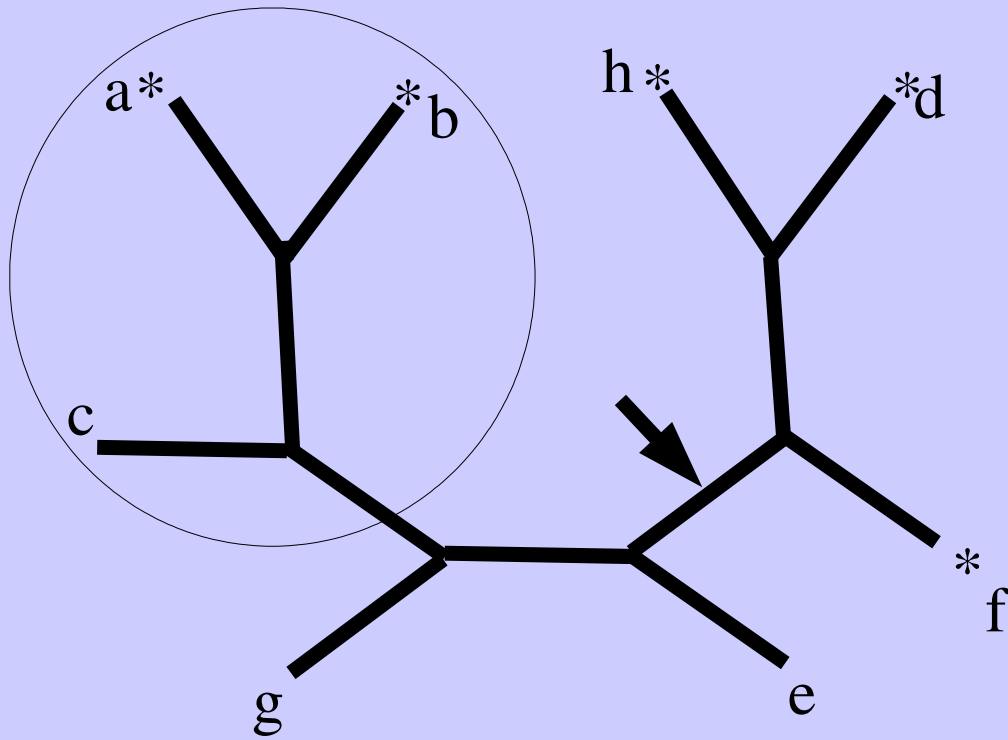Echeneis

Tetraodontiformes

Pleuronectiformes

# Taking into account probabilities of clades to appear by chance

- One must measure them for a given tree reconstruction program:

  - produce data without signal (randomly chosen character states)

  - analyze it and count the clades

  - reproduce the experiment many times
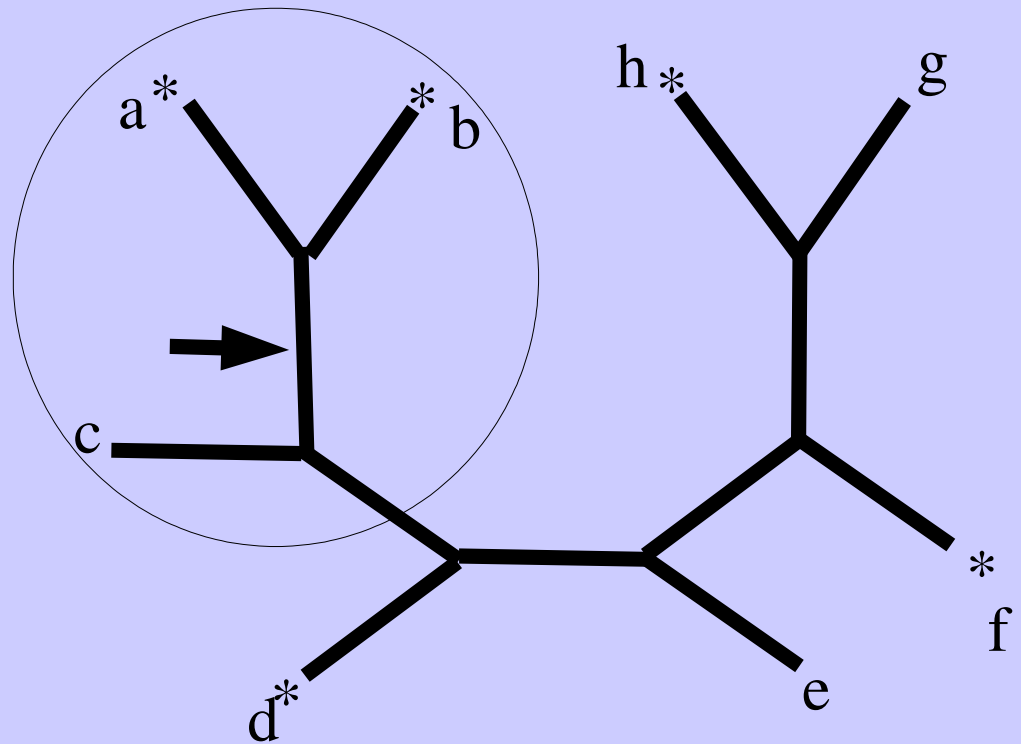
- Problem: it takes a lot of time

# Estimating probabilities of clades knowing their sizes

- All groups of the same size are equivalent

  - random data

  - random addition sequence

  - work on unrooted trees

- Shift to rooted trees

  - clades that appear in rooted trees were there before rooting

  - but clades present in unrooted trees may disappear at rooting if there are 4 outgroups or more -> one overestimates probabilities of clades containing outgroup taxa

* : outgroup taxon

a*     *b     h*     *d

c

g       e       *f

The clade is in the unrooted
tree and in the rooted one

a*     *b     h*     g

c

d*       e       *f

The clade is in the unrooted
tree but not in the rooted one

taxon 1 ATCCGTGGCAATCGG...
taxon 2 CCTAGGTGGCGAAAT...
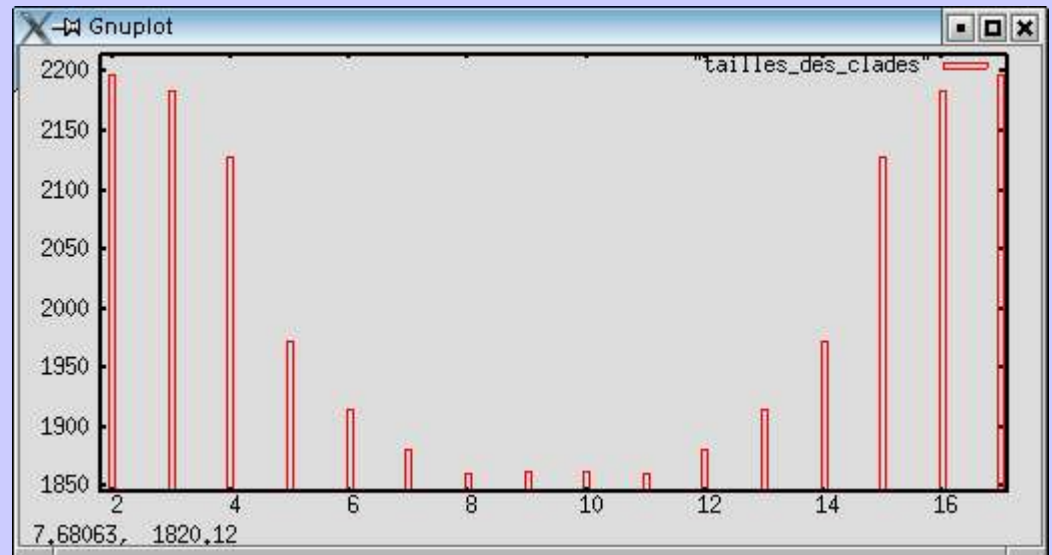taxon 3 AACATTGCGAACCTC...
...

**Phylogenetic analysis program**

*This is a black box*

Frequencies of the clades:

```
1  1.0
2  0.0128452241715
3  0.00225390436877
4  0.000548933608531
5  0.000169633642931
6  7.05808639245e-05
7  3.73445397052e-05
8  2.46185158724e-05
9  2.01545100926e-05
10 2.01545100926e-05
11 2.46185158724e-05
12 3.73445397052e-05
13 7.05808639245e-05
14 0.000169633642931
15 0.000548933608531
16 0.00225390436877
17 0.0128452241715
18 1.0
19 1.0
```

Gnuplot

"tailles_des_clades"

2200
2150
2100
2050
2000
1950
1900
1850

2    4    6    8    10    12    14    16

7,68063,  1820,12

# From counts to probabilities

N(t): number of clades including t taxa

T: number of trees examined

N(t)/T = mean number of clades including t taxa per tree

n: number of taxa in the trees

n!/t!(n-t)!: number of possible clades including t taxa when there are n taxa in the study

P(t) = (N(t)/T)/(n!/t!(n-t)!)