

# *Evolution des Eucaryotes*

## L'éclairage de la phylogénie moléculaire



Denis BAURAIN

Unité de Phylogénomique des Eucaryotes  
Département des Sciences de la Vie / ULg

Collège Belgique / Palais des Académies  
23 octobre 2014



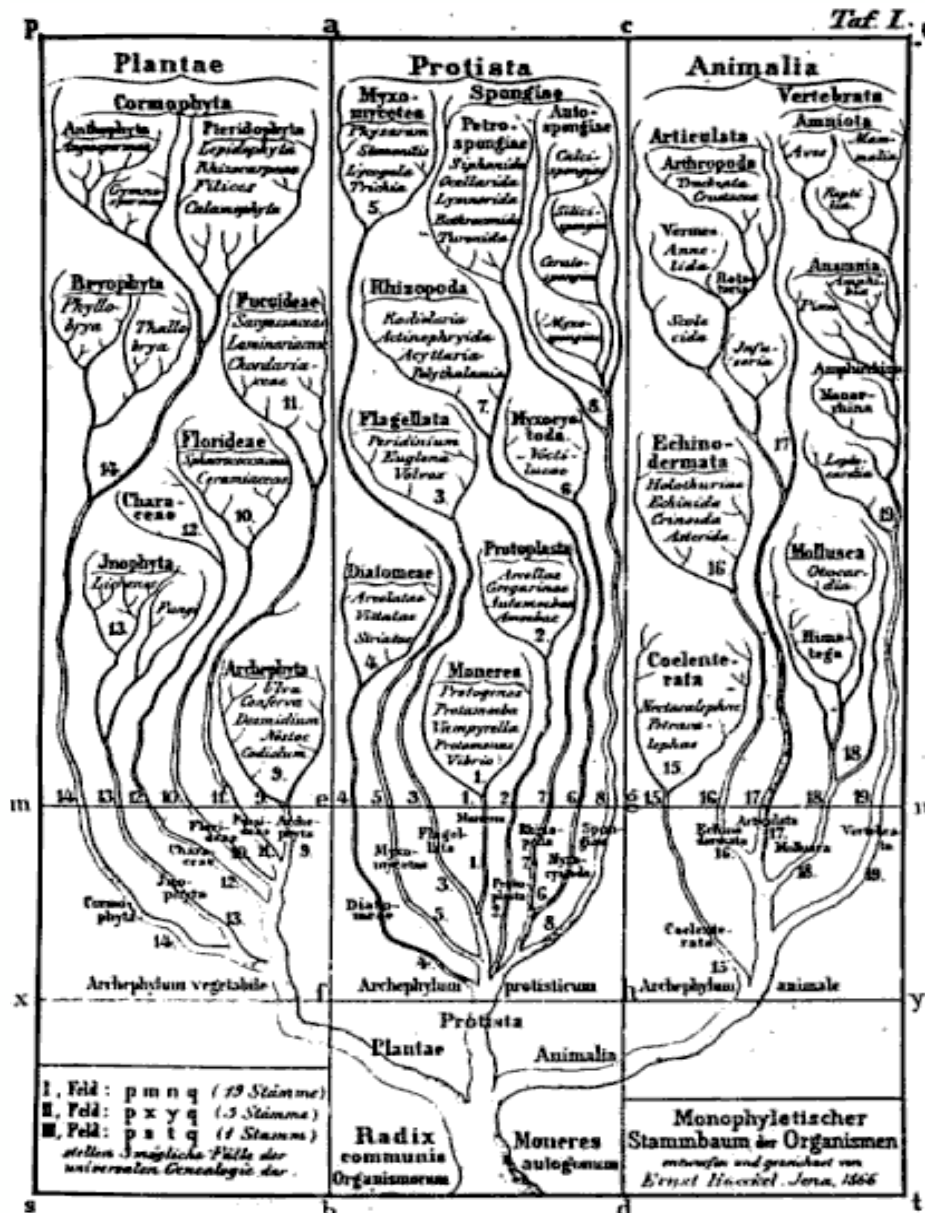
*Phylogénie moléculaire  
Grandeur et décadence  
des Archezoa*

# Phylogénie



Ch. Darwin

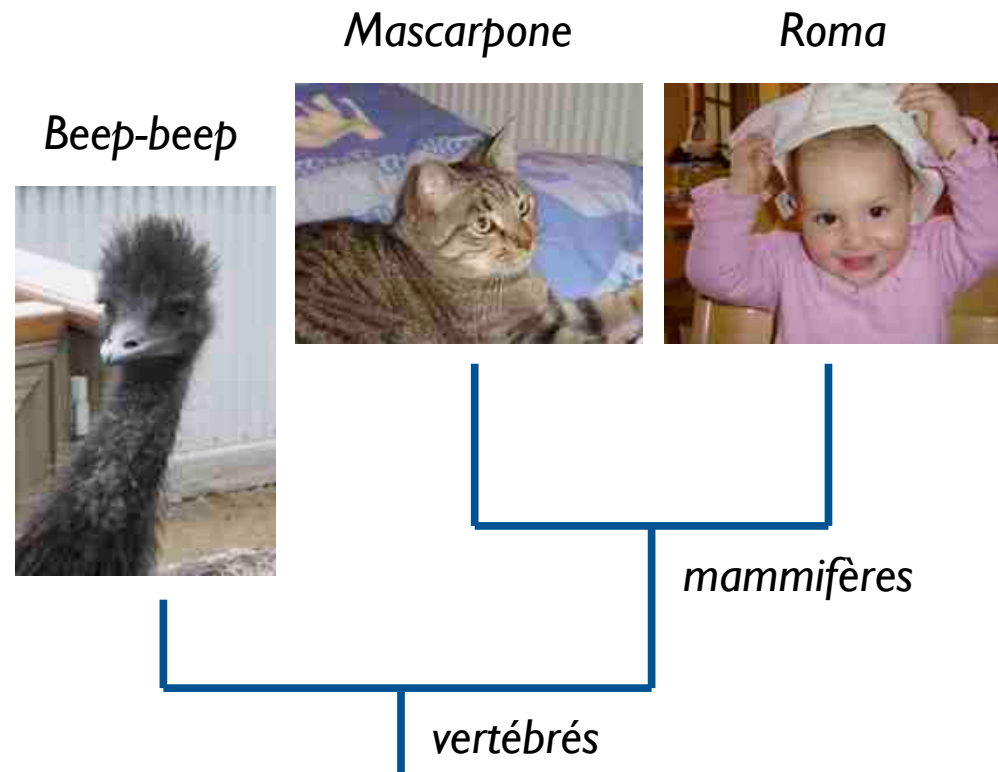
E. Haeckel



- domaine de la biologie étudiant les **relations de parenté** entre organismes actuels
- fondée sur la théorie de **l'évolution des espèces** de Ch. Darwin (1859)
- utilisant **l'arbre** à la fois comme métaphore (E. Haeckel, 1866) et outil de représentation

# Phylogénie classique

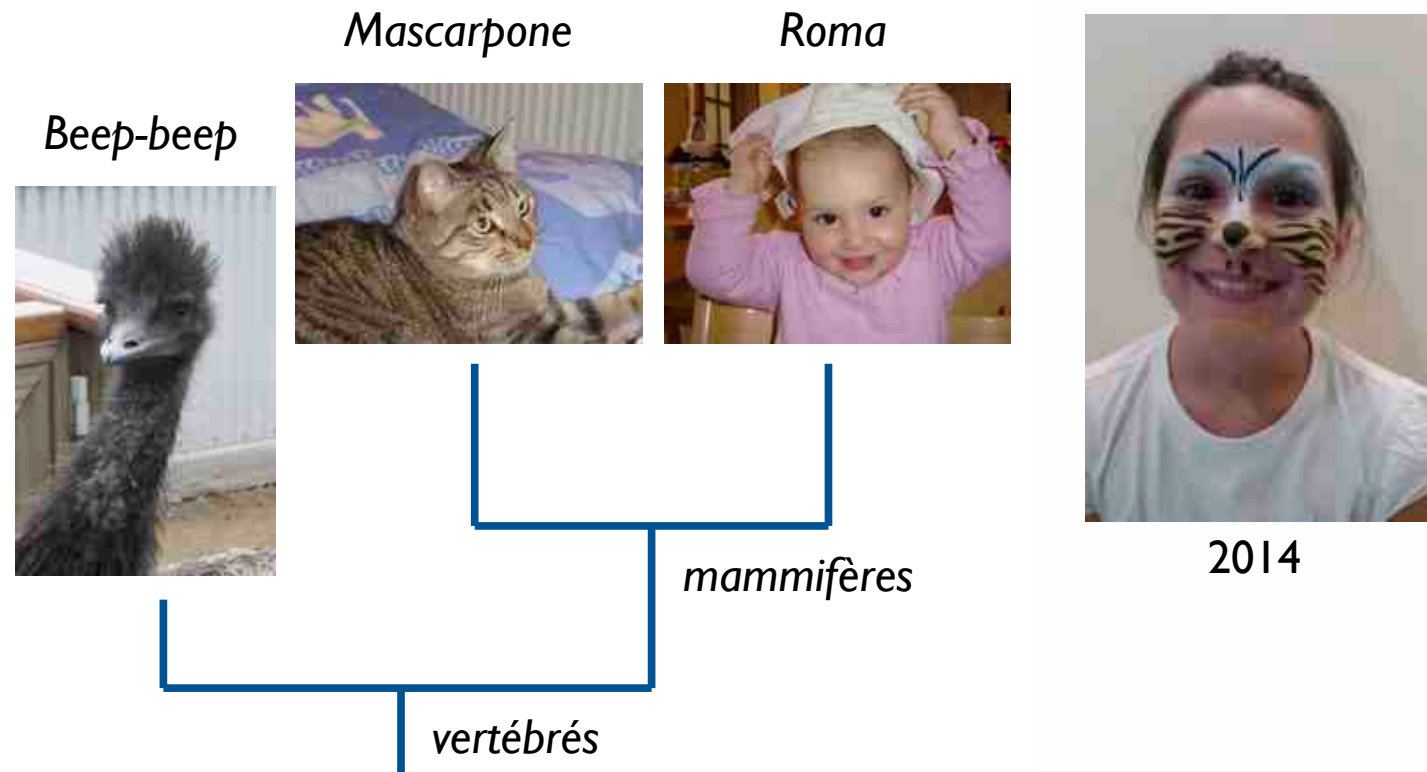
- Caractères morphologiques
  - ouverts aux interprétations
  - pas toujours disponibles (ex. microbes)
  - souvent incommensurables entre lignées





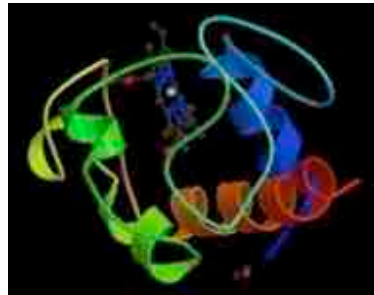
# Phylogénie classique

- Caractères morphologiques
  - ouverts aux interprétations
  - pas toujours disponibles (ex. microbes)
  - souvent incommensurables entre lignées

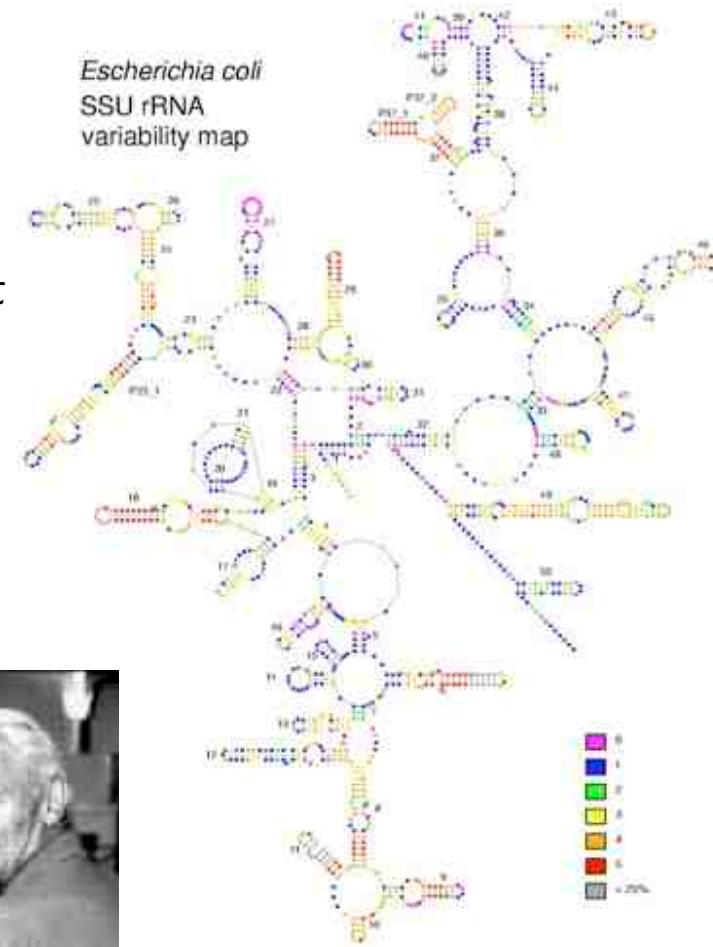


# Phylogénie moléculaire

- Caractères moléculaires
  - protéines
    - Fitch & Margoliash (1967)



cytochrome c  
de cheval



- DNA
- SSU rRNA (16S/18S)
  - Carl Woese (1977)
  - molécule essentielle
  - structure universelle
  - très ancien



Carl Woese



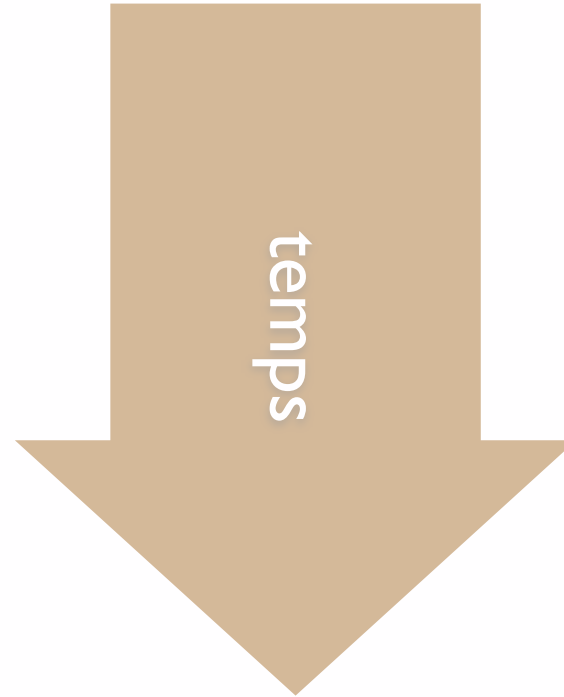
Walter M. Fitch

# Phylogénie moléculaire

principe de base / raisonnement de parcimonie

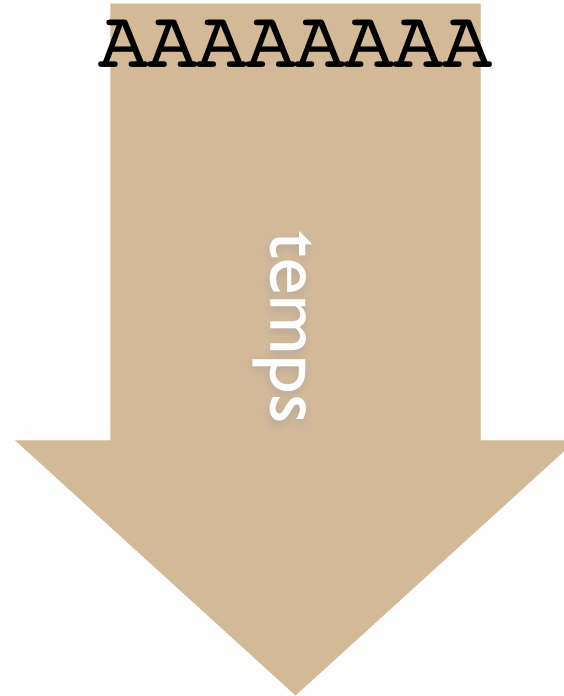
# Phylogénie moléculaire

principe de base / raisonnement de parcimonie



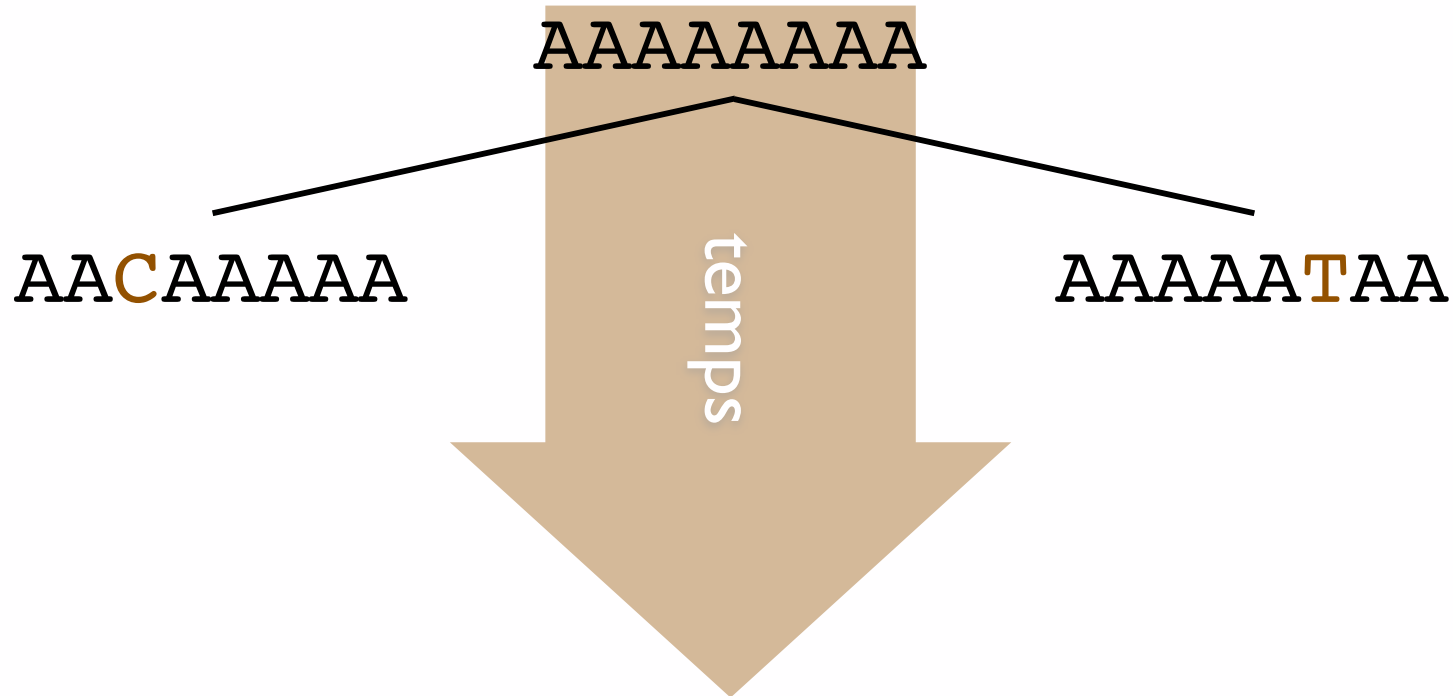
# Phylogénie moléculaire

principe de base / raisonnement de parcimonie



# Phylogénie moléculaire

principe de base / raisonnement de parcimonie

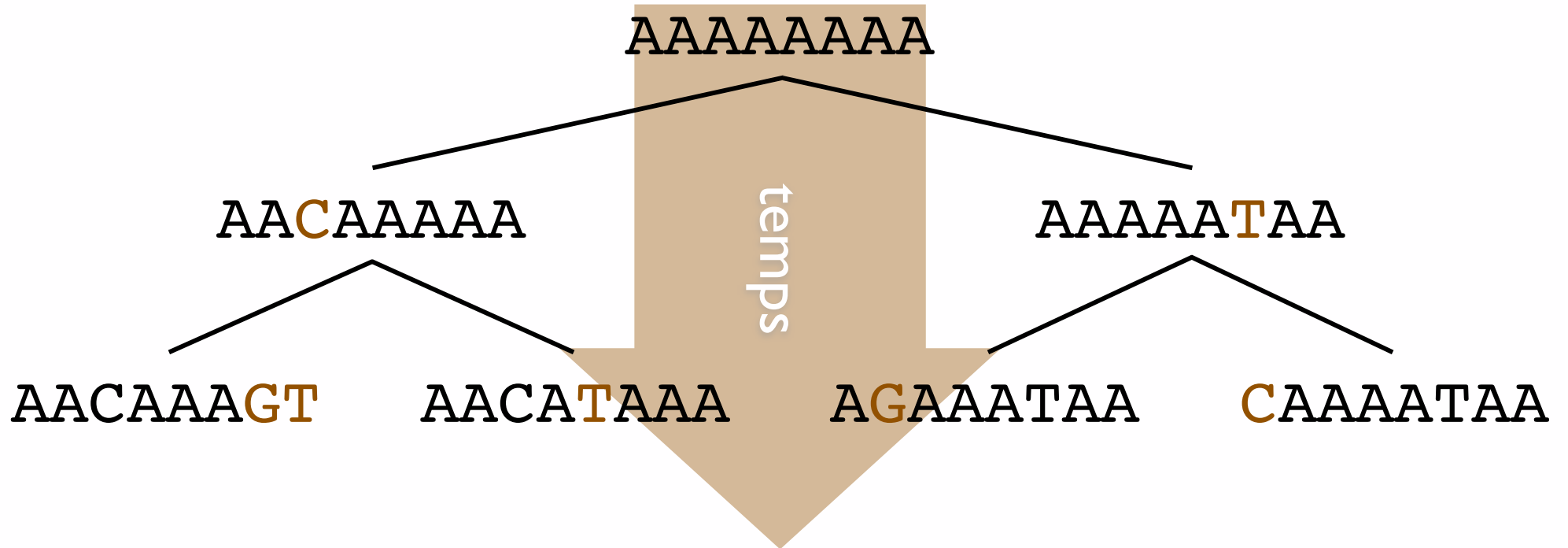






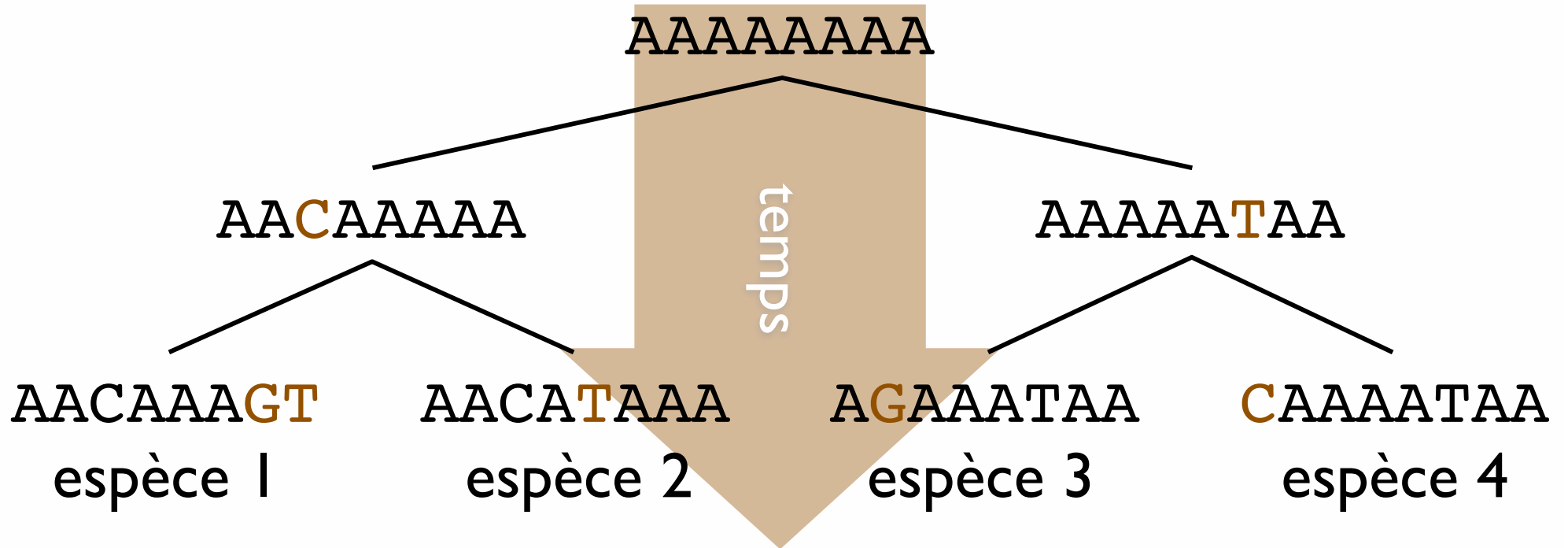
# Phylogénie moléculaire

principe de base / raisonnement de parcimonie



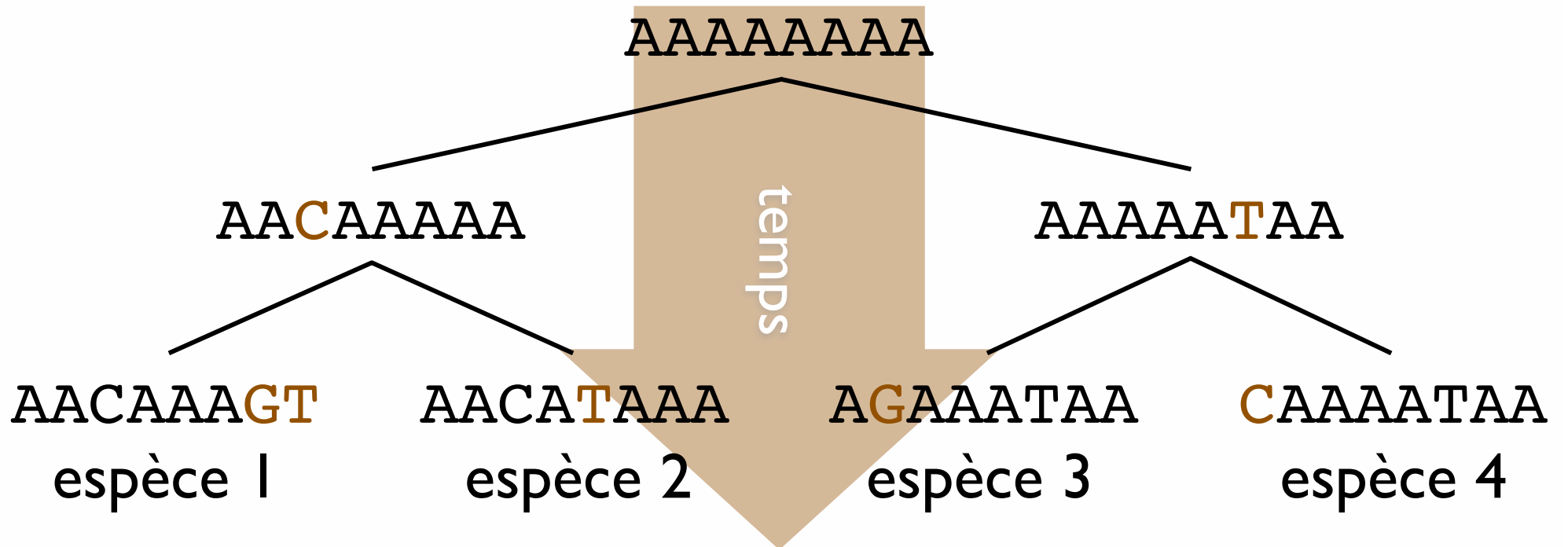
# Phylogénie moléculaire

principe de base / raisonnement de parcimonie



# Phylogénie moléculaire

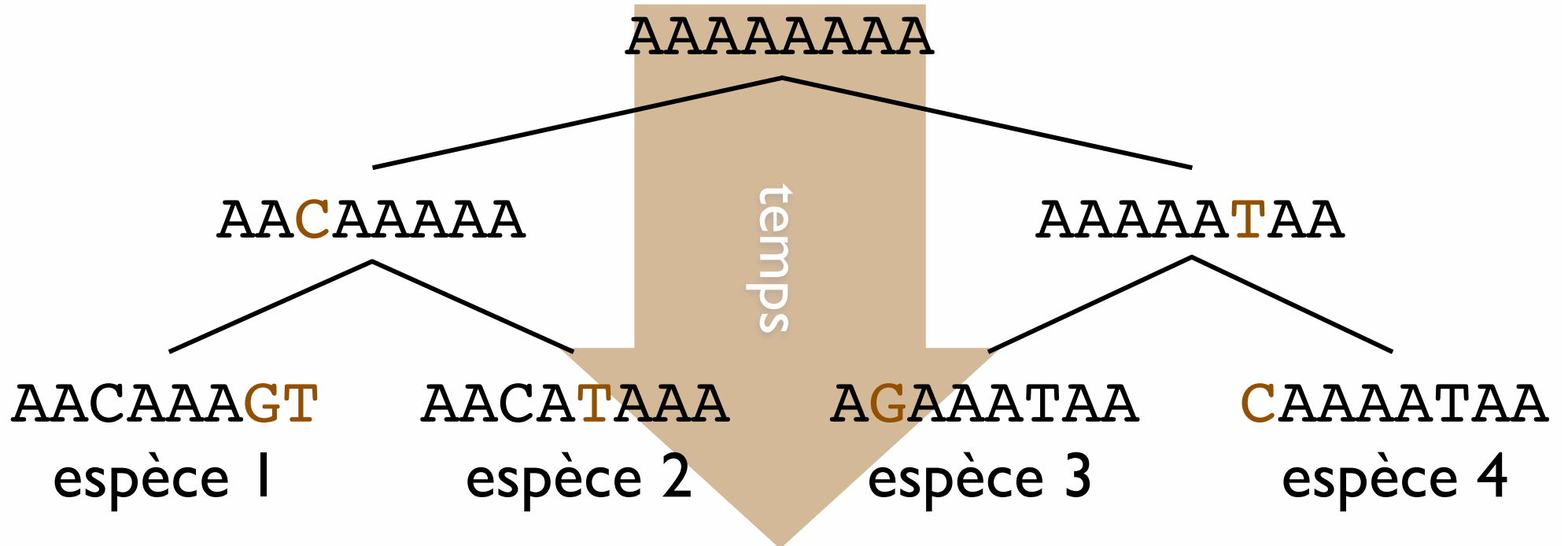
principe de base / raisonnement de parcimonie



- 1 AAC**A**AAAGT
- 2 AAC**A**TAAA
- 3 AGAAA**T**AA
- 4 CAAA**A**TAA

# Phylogénie moléculaire

principe de base / raisonnement de parcimonie

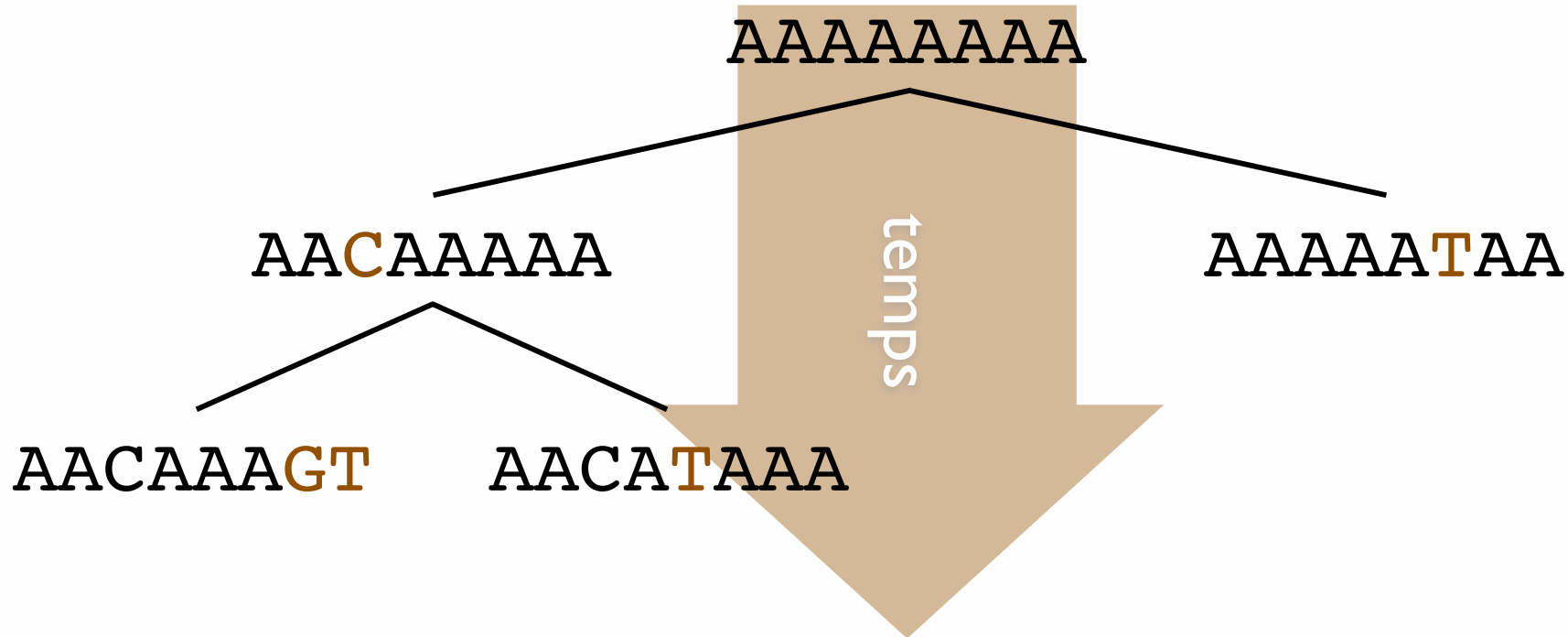


- 1 AACAAAGT
- 2 AACATAAA
- 3 AGAAATAA
- 4 CAAAATAA

Les substitutions sont assez **rare**s.  
La majorité des différences de séquence observées entre espèces sont donc **héritées**.

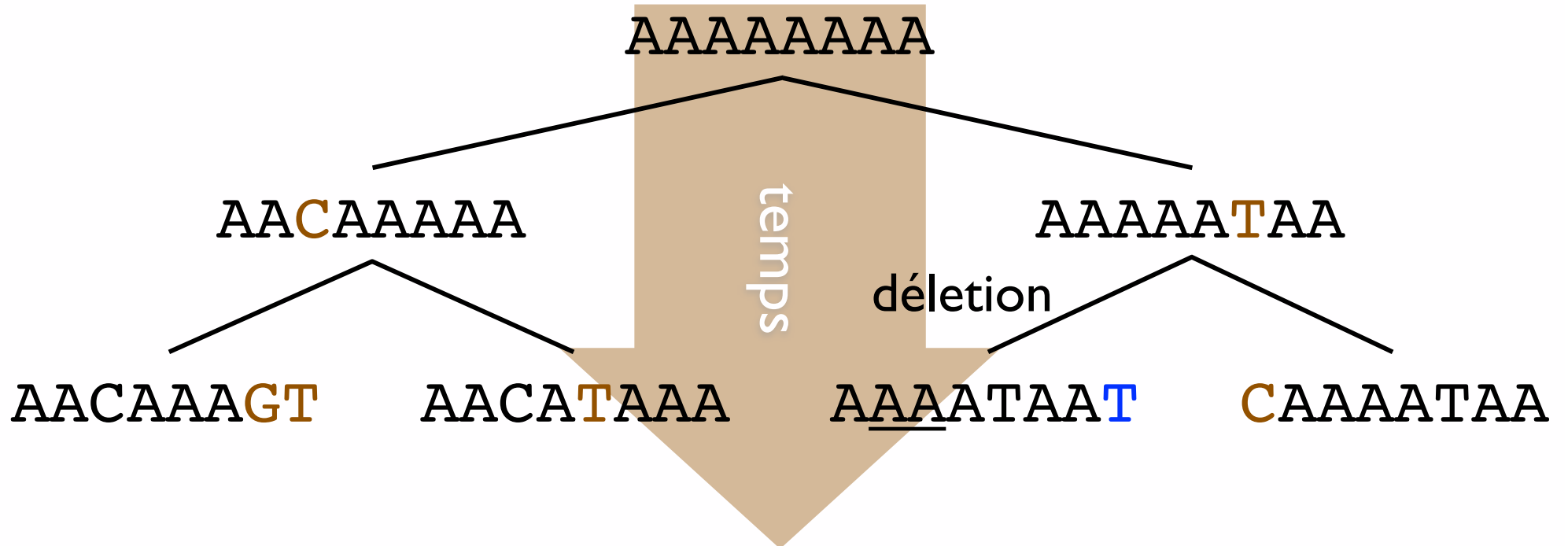
# Phylogénie moléculaire

nécessité d'aligner les séquences



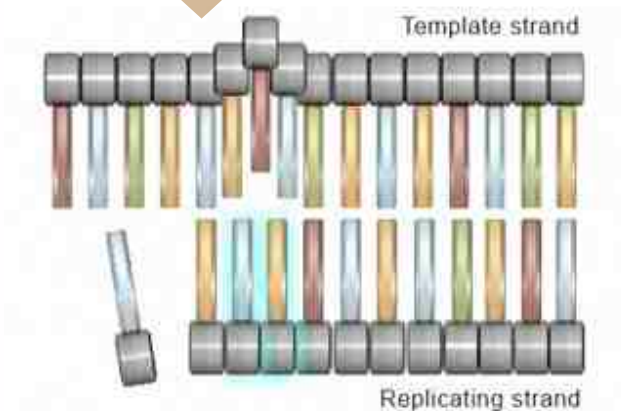
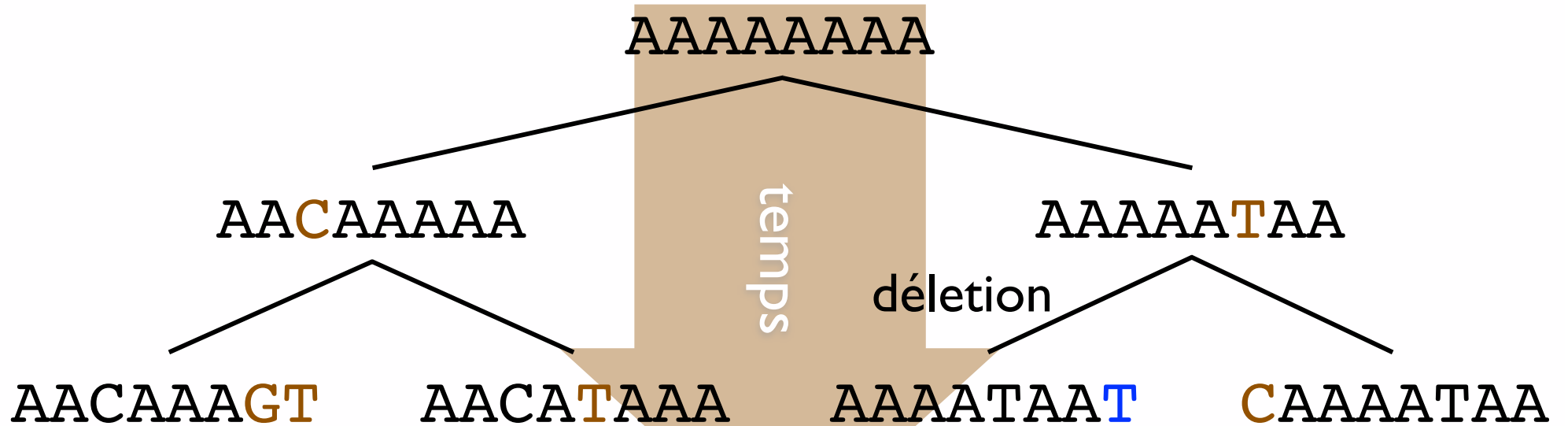
# Phylogénie moléculaire

nécessité d'aligner les séquences



# Phylogénie moléculaire

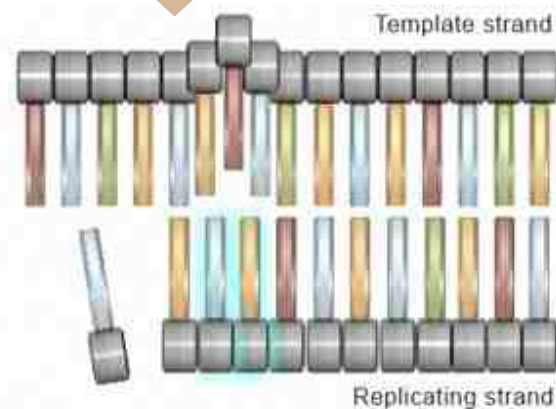
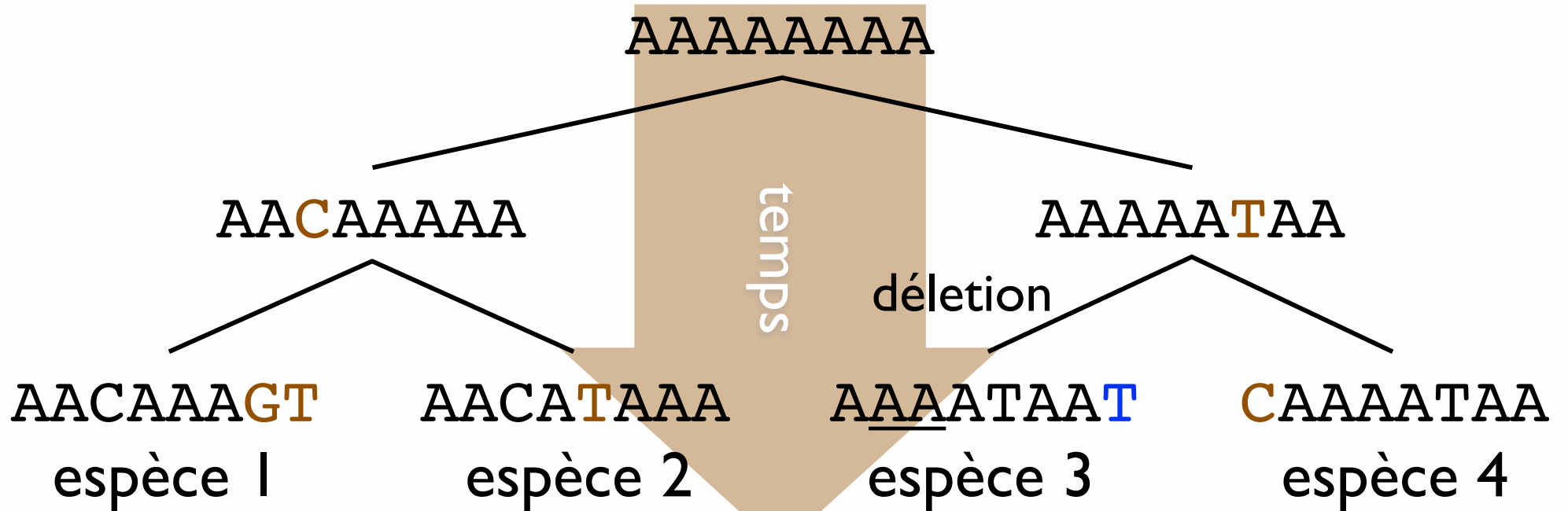
nécessité d'aligner les séquences





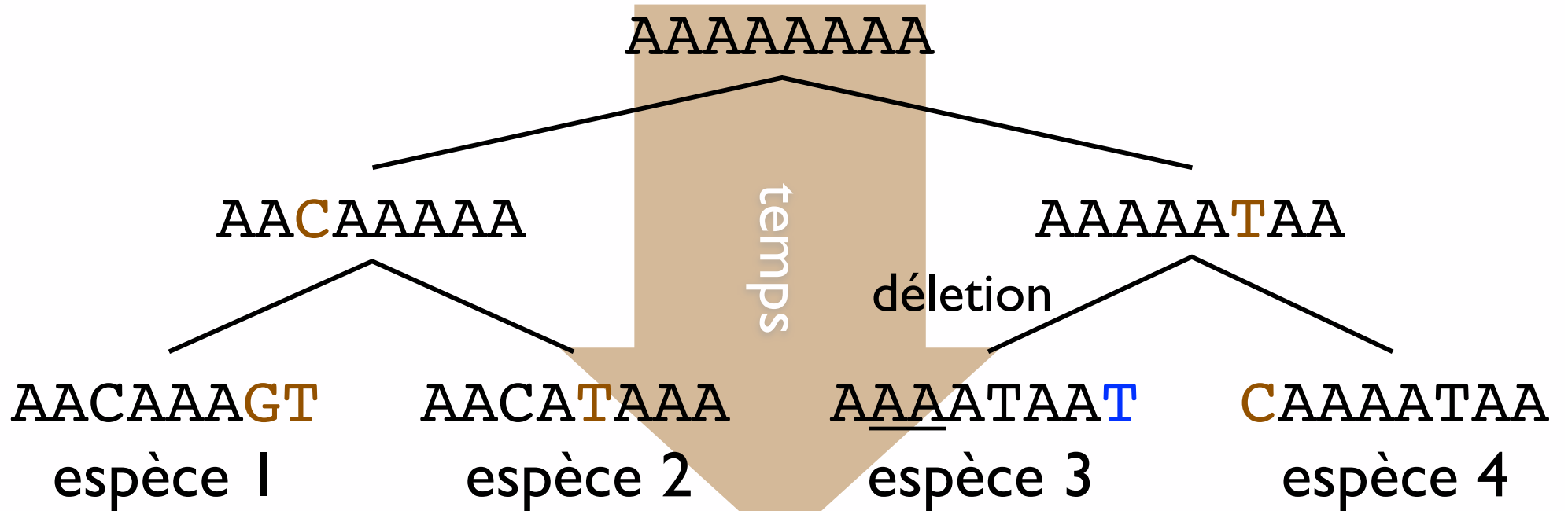
# Phylogénie moléculaire

nécessité d'aligner les séquences

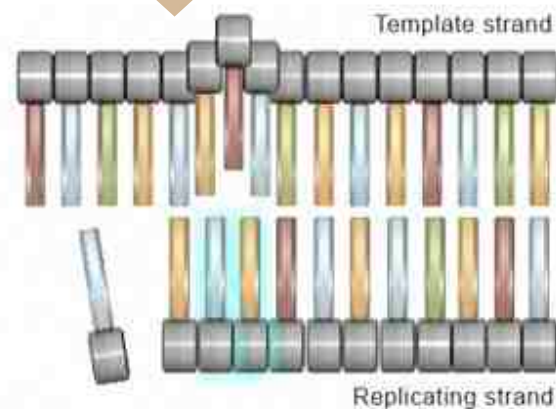


# Phylogénie moléculaire

nécessité d'aligner les séquences

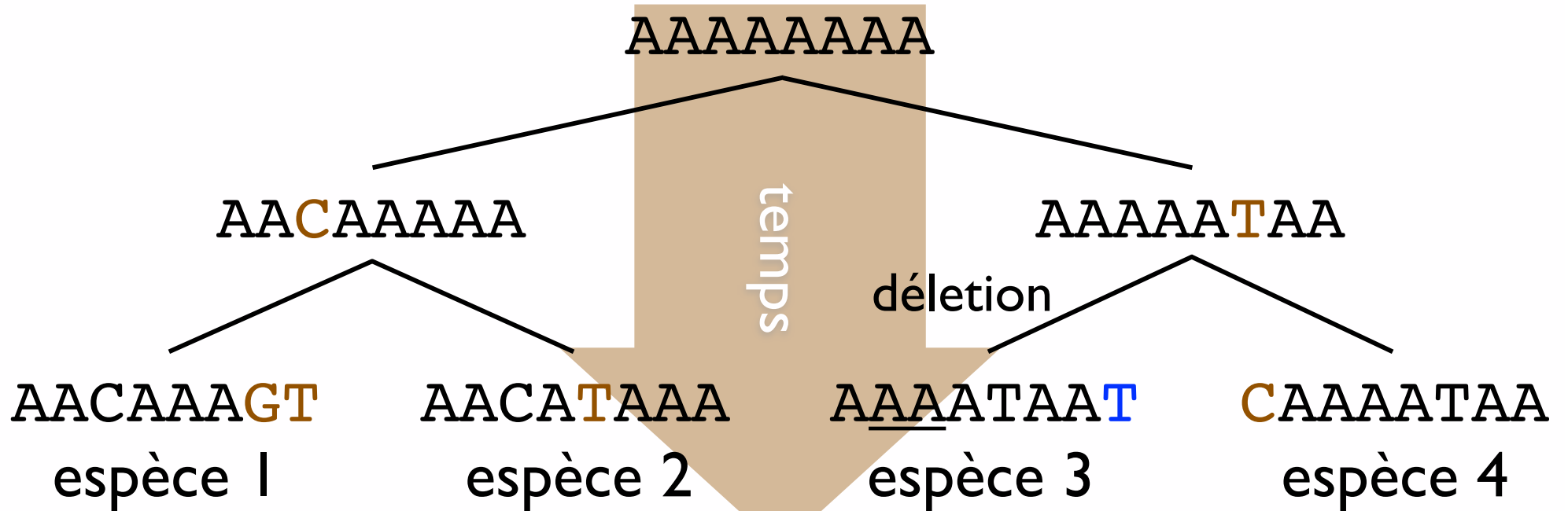


- 1 AACAAAGT
- 2 AACATAAA
- 3 AAAATAAT
- 4 CAAATAA

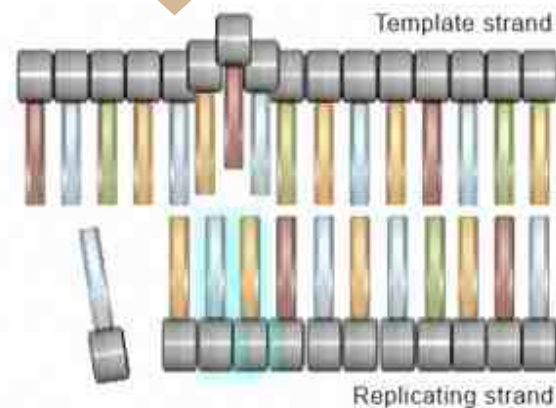


# Phylogénie moléculaire

nécessité d'aligner les séquences



1 AACAAAGT  
2 AACATAAA  
3 AAAATAAT  
4 CAAAAATAA



1 AACAAAGT  
2 AACATAAA  
3 AA-AAATAA  
4 CAAAAATAA

# Phylogénie moléculaire

## exemple : oxydase alternative

```

          100      110      120      130      140      150      160      170      180
Cr1  FAAGAVA P HPGINPARMAADSASAAAGASGDAALAESYMAHPAYSDEYVESVRPTHTVTPOKLHOHVGLR T IOVFR YLFDKA TGY TPTGS-
Cr2  FATSGVA P HPGMKAPSPPTDDEVEACW-----RPVYDTAYLEK V KPFHIT P ERLYORIGFRAIMAARWTFDKL TGY GP--N-
Ac   AATSNSNMRYF S STSRRWIKEFFA P PKET D HIVESVTTWKHPVFT EKQOMKEIAIAHREAKNWS DWV ALGTV RFLRWATDLA TGY RHAA--
An   NRHQTAGKRFI S TTPKSOIKEFF- P PPTAPHVKEVETAWVHPVY TEEQMKQVAIAHRD AKNWADWV ALGTV RMLRWGMDLV TGY RHPP--
Ca   EKPGTIPTKHKP F NIQTEVYNKAGIEAND D DKFLT KPTYR HEDF TEAGVYR VHVTHRPPRTIGDKI SCYGT LFFR KCFDLV TGY AVDPD-
Mg   LPRLAAS P RLF S TTSSAQLRDF- P VKETEHIROTPPTWPHHGL TEKEMVDV VPGHRK PRTLGD KFAWSL V RISRWGMDKV SGLSSEQQ
Nc   TNLSSPS P RNF S TTSVTRLKDF- P AKETAYIROTPPAWPHHGW TEEEMTS V VPEHRK P ETVGDWL AWKL V RICRWATDIA TGI RPE-QQ
Pa   PSPHSKD P NSK S IFDIGTKLIVNP P PQMAD NQYVTHPLFPHPKY SDEDCEAV H FVHRE P KTI GDKIADR GV KFCRAS FDFV TGY K KPKDV
At   GGDAAGGNNKGD KGIAS YWGVEPNKITKEDGSEWKWNCFRPWET YKADITIDLK KHHV P TTF LDR IAYWTV KSLRWPTDL-----
Sg   DGGAE-----KEAVV SYWAVPPSKVSKEDGSEWRWTCFRPWET YQADLSIDLHKHHV P TTI LDKL ALRTV KALRWPTDI-----

```

\*

```

          190      200      210      220      230      240      250      260      270
Cr1  -----MTEAQWL RRM I FLET VAGCPGMVAGMLRHL KSLRSM SRD R GWI HTLL EEAENERMHL I TFLQLRQP GPAFRAMVI
Cr2  -----MTEAKWL ORM I FLET IAGVPGMVAGVLRHL KSLRSM KR D H GWI HTLL OEAENERMHL L TFFELRKP GP LFRASII
Ac   --PGKQGV E VPEQFQ MTERK W V I R F I FLET VAGVPGMV G M L R H L R S L R R M K R D N G W I E T L L E E A Y N E R M H L L S F L K L A O P G W F M R L M V L
An   --PGR---EHEARFK MTEOKWL TRF I FLES VAGVPGMV G M L R H L R S L R R M K R D N G W I E T L L E E A Y N E R M H L L T F L K L A E P G W F M R L M V L
Ca   -DKPDQYKGT--RWE MTEEKW MTRC I FLES IAGVPGSVAGFVRHL HSLRML TRD KAWI ETLH DEAYNERMHL L T F I K I G K P S W F T R S I I Y
Mg   INKGSPTTSIVA AKPL T E A Q W L S R F I F L E S I A A V P G M V A G M L R H L H S L R R L K R D N G W I E T L L E E A Y N E R M H L L T F L K M C E P G W L M K I L I I
Nc   VDKHHPTTATSADKPL T E A Q W L V R F I F L E S I A G V P G M V A G M L R H L H S L R R L K R D N G W I E T L L E E S Y N E R M H L L T F M K M C E P G L L M K T L I L
Pa   NGMLKSWEGT--RYE MTEEKWL TRC I FLES VAGVPGMVA F I R H L H S L R L L K R D K A W I E T L L D E A Y N E R M H L L T F I K I G N P S W F T R F I I Y
At   -----FFORRYGCRAMMLETVA AVPGMV G M L L H C K S L R R F E O S G G W I K A L L E E A E N E R M H L M T F M E V A K P K W Y E R A L V I
Sg   -----FFORRYACRAMMLETVA AVPGMV G V L L H L K S L R R F E H S G G W I R A L L E E A E N E R M H L M T F M E V A Q P R W Y E R A L V L

```

\*

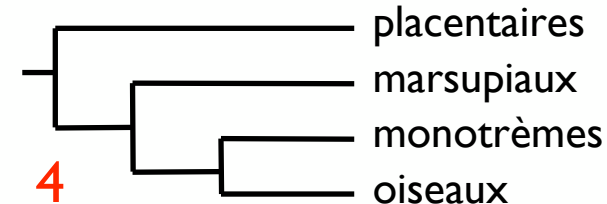
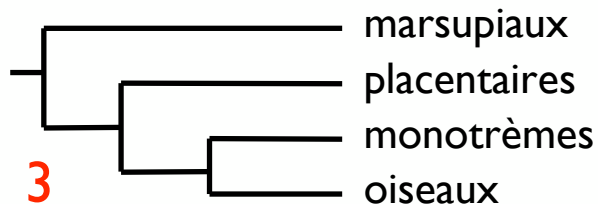
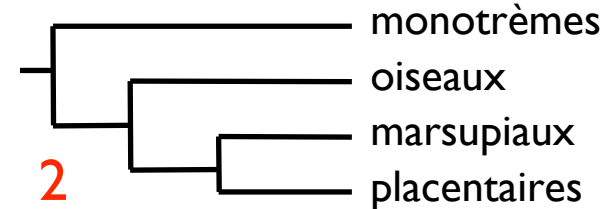
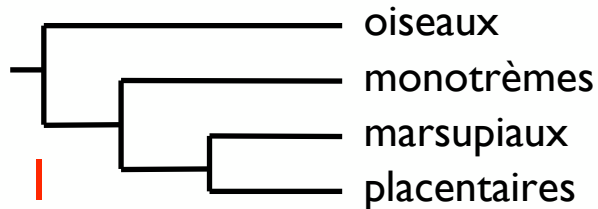
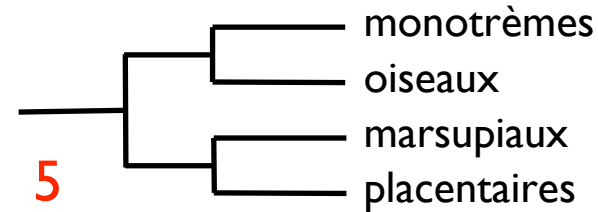
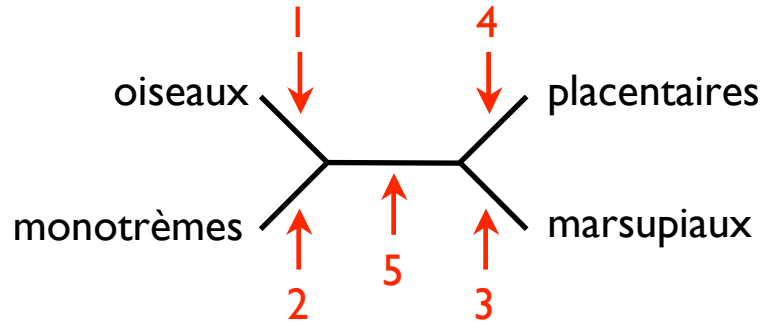
```

          280      290      300      310      320      330      340      350      360
Cr1  LAQGVFFNAYF IAYLLSPRT CHAFVGFLEEEAVK TYTHALVEIDAGRL----WKDTPAPPVAVOYWGLKPG-ANMRDLIILAVRADEACHA
Cr2  VAQGVFWNLYF IGYLVSPRT CHAAVGFLEEEAVK TYTHALQEIDAGRL----WKGKVAPPIACEYWGLKPG-ASMRDLIILAVRADEACHA
Ac   GAQGVFFNGFFI SYLISPR TCHR FVGYLEEEAVM TYTHAIKDLES GKLPNW--ANQPAPDIAVAYWQMP E GKRTI L D L L Y I R A D E A K H R
An   GAQGVFFNGFFL SYLMSPRI CHR FVGYLEEEAVI TYTRAIKEIEAGSLPAW--EKTEAPEIAVOYWKMPE GORSMKD L L Y I R A D E A K H R
Ca   IGQGVFTNIFFL VYLMNPRY CHR FVGYLEEEAVR TYTHL I D E L D D P N K L P - D F O K L P I P N I A V O Y W P E L T P E S S F K D L I L R I R A D E A K H R
Mg   GAQGVYFNAMFVAYLISPKI CHR FVGYLEEEAVHTYTRSTEELERGDLPKWSDPKFOVPEIAVS YWGMPEGHRTMTRD L L Y I R A D E A N H R
Nc   GAQGVFFNAMFL SYLISPKI THRFVGYLEEEAVHTYTRCIREIEEGHLPKWSDEKFEIPEMAVRYWRMPE GKRTM K D L I H Y I R A D E A V H R
Pa   MGQGVFANLFFLVYLIKPRY CHR FVGYLEEEAVS TYTHL I K D I D S - K R L P - K F D D V N L P E I S W L Y W T D L N E K S T E R D L I Q I R A D E S K H R
At   TVQGVFFNAYFLGYLISPKFAHRMVGYLEEEA I H S Y T E F L K E L D K G N I ----ENV PAPAIAIDYWRLPAD-ATL R D V V M V V R A D E A H H R
Sg   AVQGVFFNAYFLGYLISPKFAHRV VGYLEEEA I H S Y T E F L K D I D N G A I ----QDC PAPAIALDYWRLPOG-STL R D V V T V V R A D E A H H R

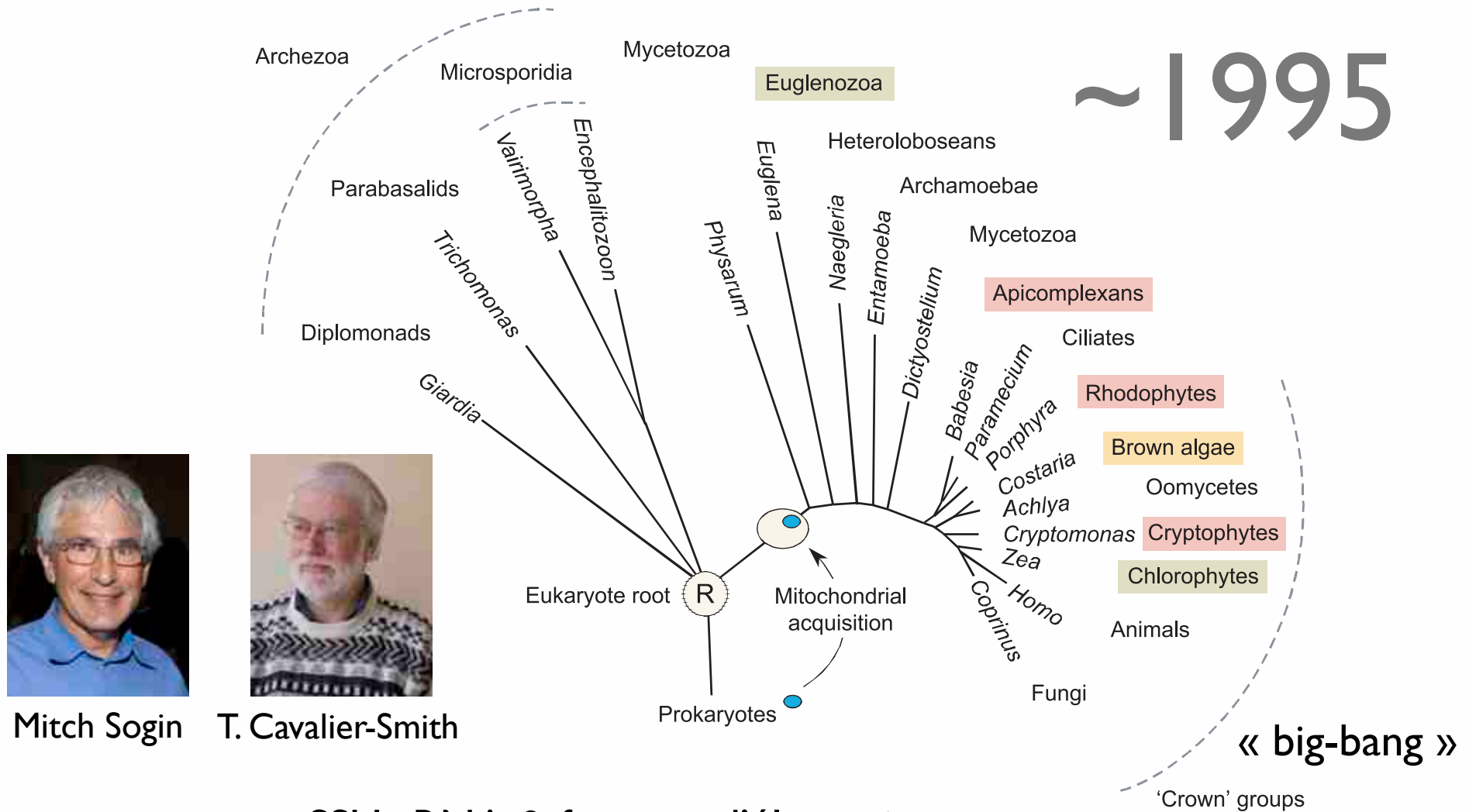
```

# Phylogénie moléculaire

arbre non-enraciné, racine et groupe extérieur



# Archezoa vs. 'Crown'



Mitch Sogin



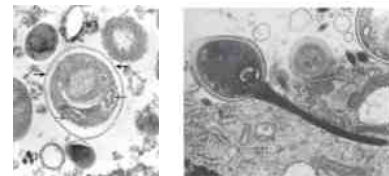
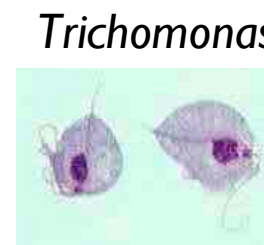
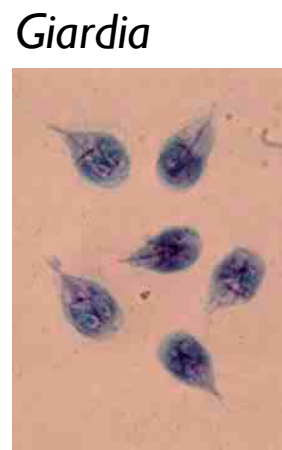
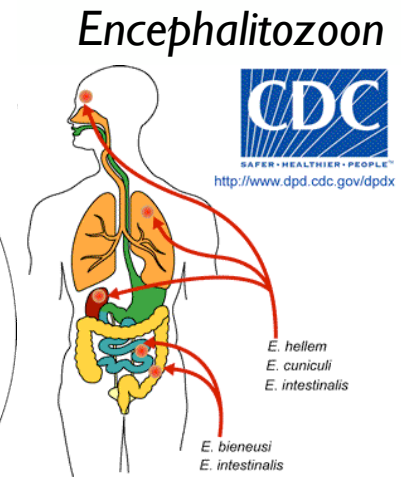
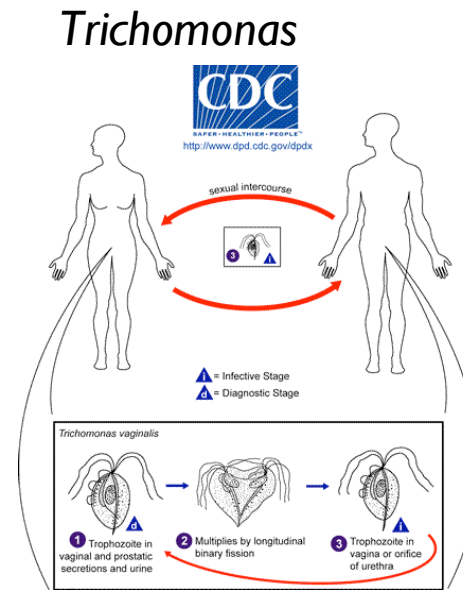
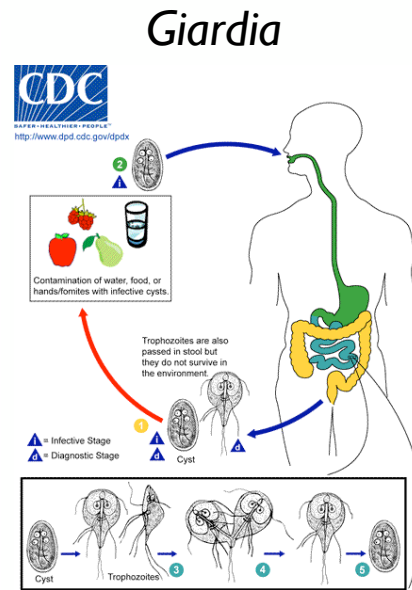
T. Cavalier-Smith

## SSU rRNA & facteurs d'élongation

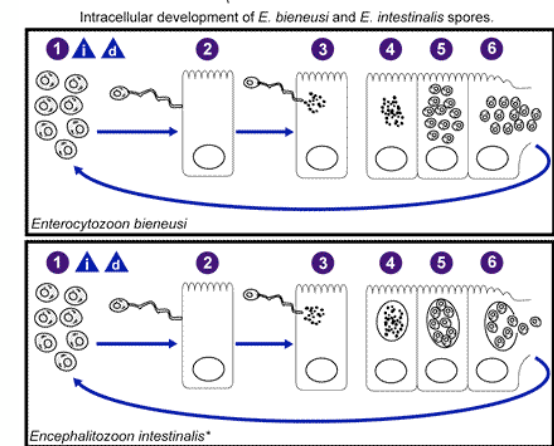
# Que sont les Archezoa ?

parasites  
amitochondriaux

(ultrastructure simple,  
pas de peroxyosomes,  
système endomembranaire  
peu développé)



Encephalitozoon



\*Development inside parasitophorous vacuole also occurs in *E. hellem* and *E. cuniculi*.



# Pas si Archezoa (I)

[incongruence phylogénétique]

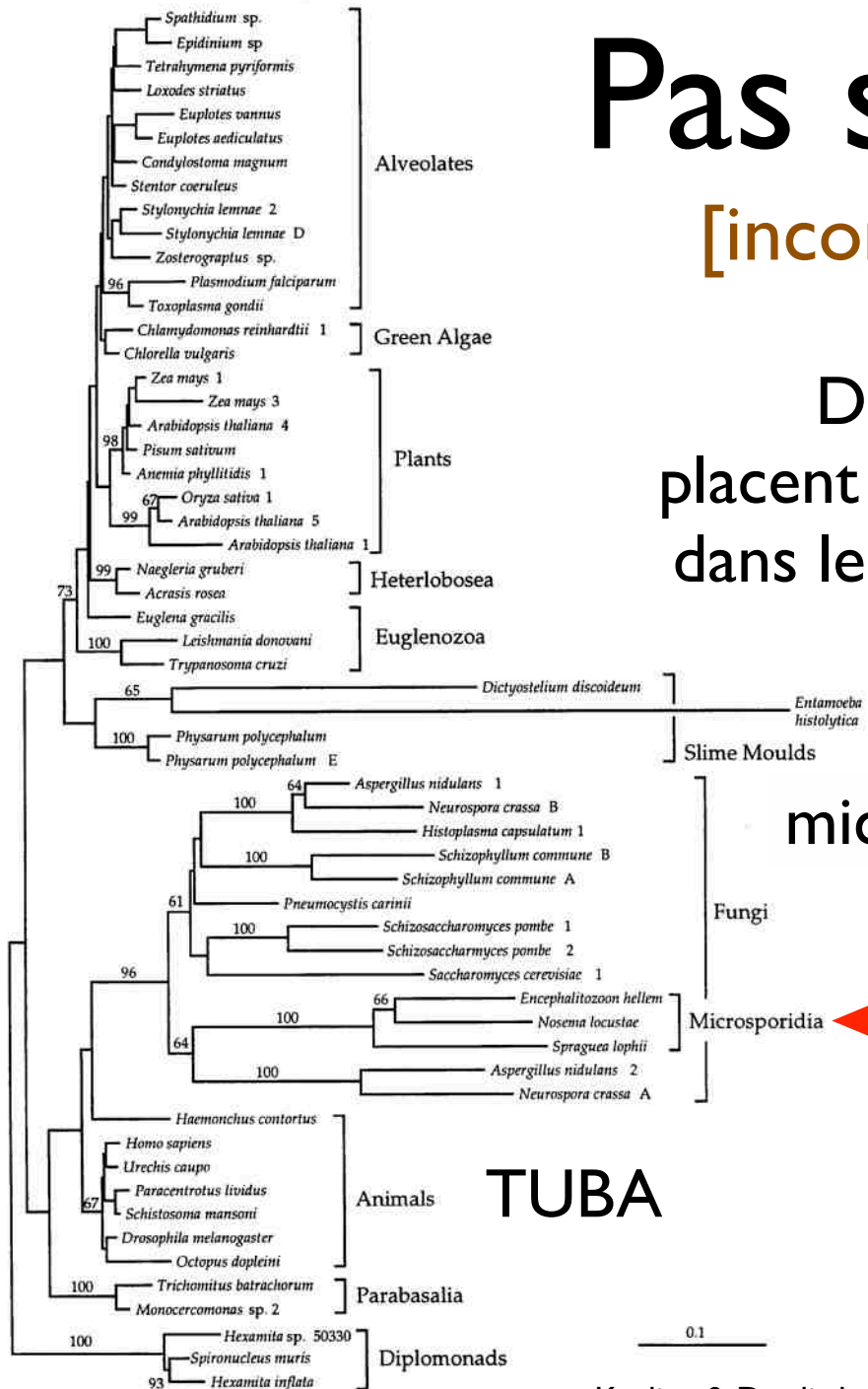
D'autres gènes placent les Archezoa dans le **crown group**.



Patrick Keeling



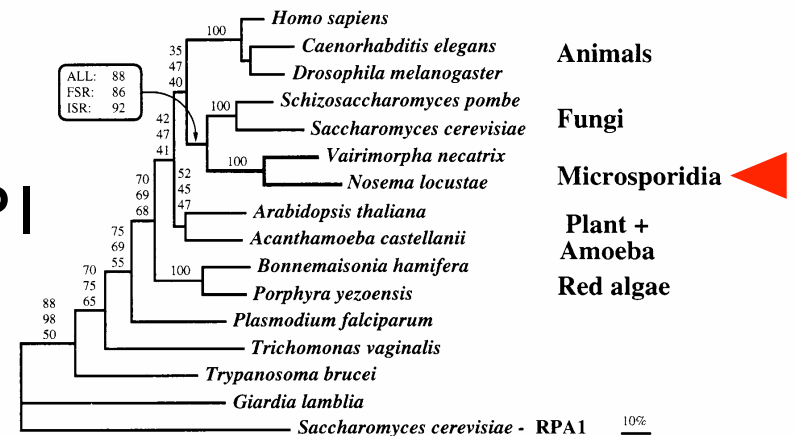
Martin Embley



microsporidies = champignons dérivés

RBPI

TUBA





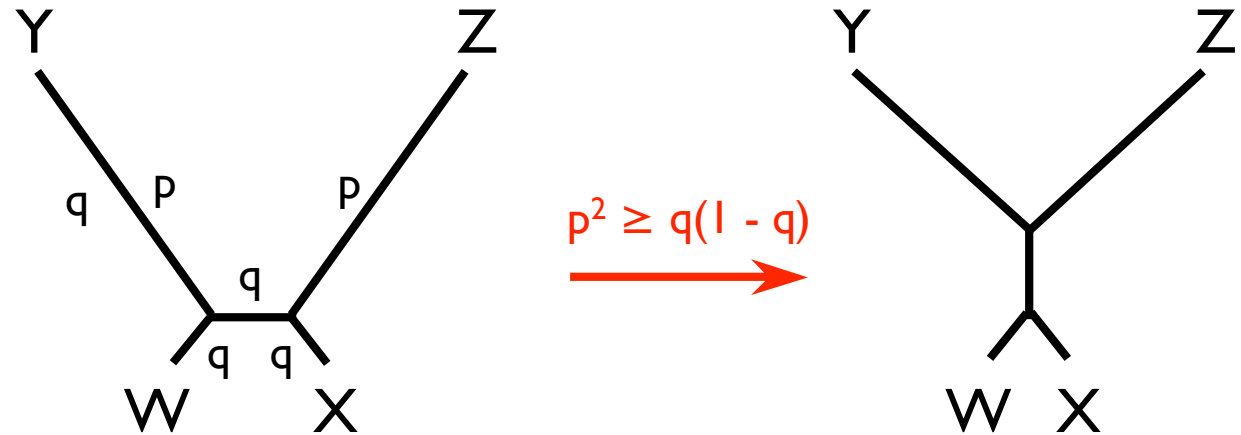
Joe Felsenstein



Hervé Philippe

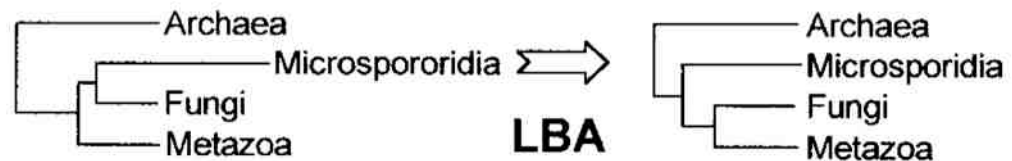
# Pas si Archezoa (2)

[artéfacts de reconstruction]



si Y est un outgroup phylogénétiquement éloigné

L'attraction des longues branches (LBA) est responsable de la **position artéfactuelle** des Archezoa dans l'arbre des eucaryotes.





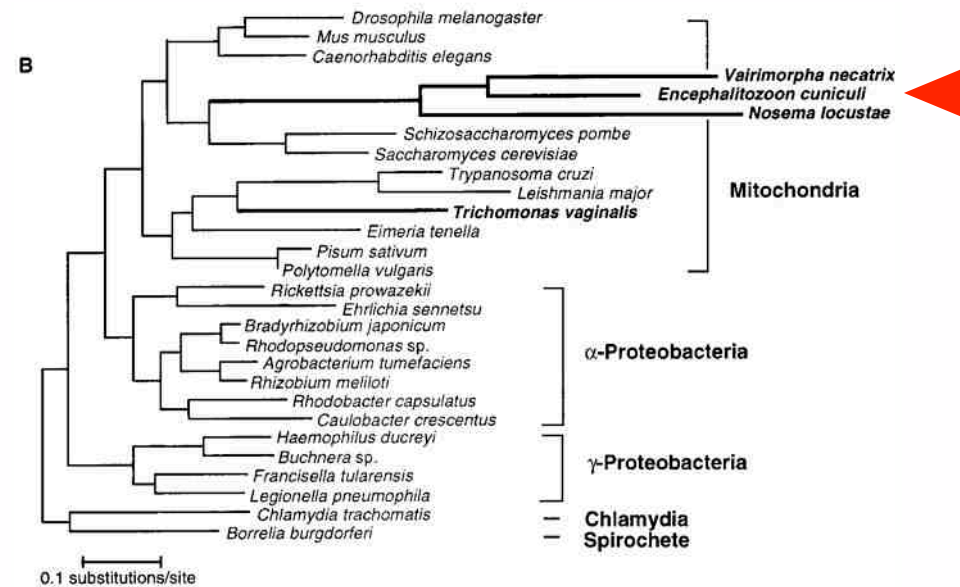
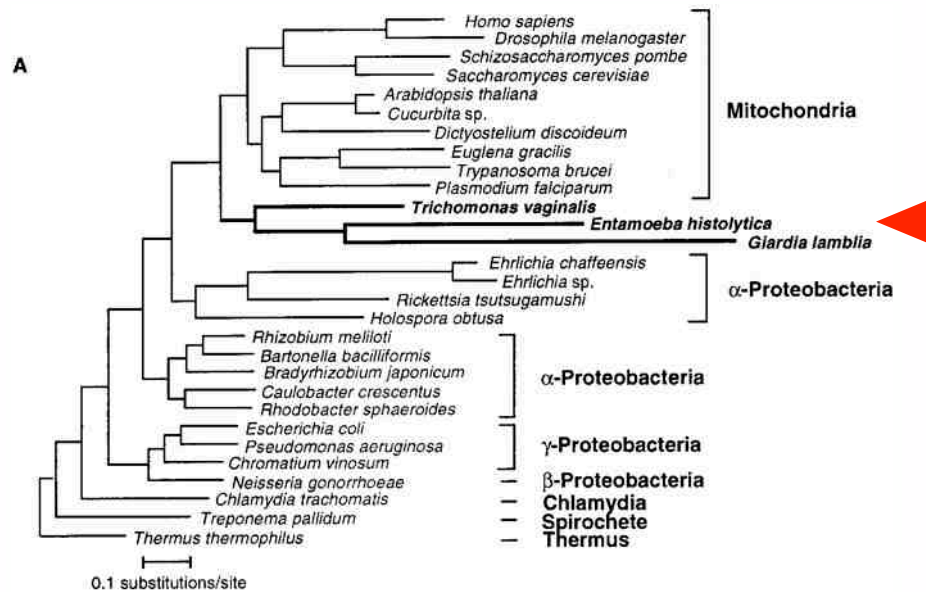
Andrew Roger

# Pas si Archezoa (3)

[transfert de gènes]

## CPN60

## HSP70



présence de gènes nucléaires d'origine mitochondriale dans les Archezoa

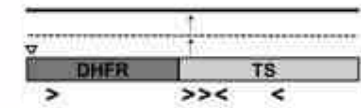


# Plan

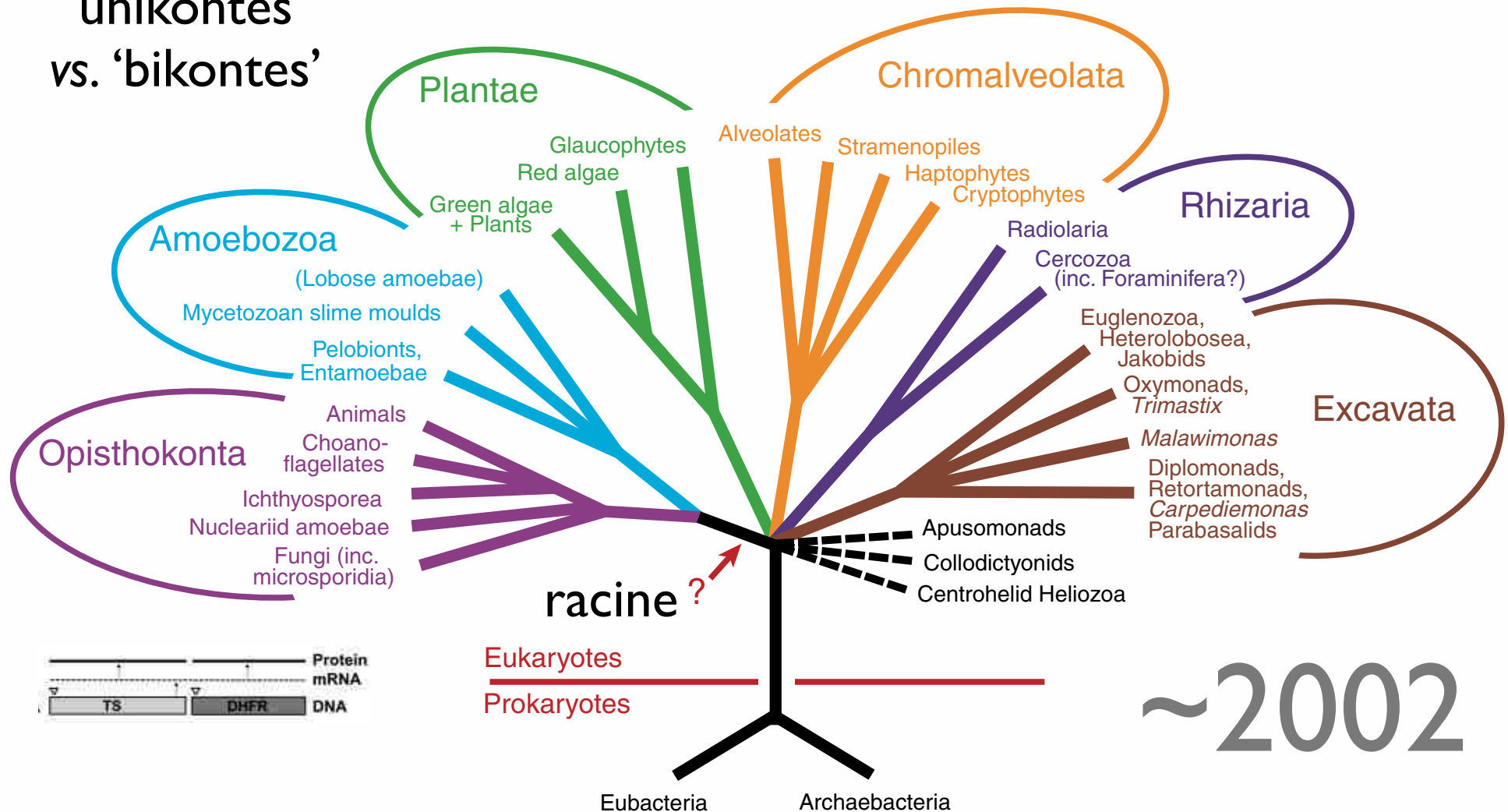
1. Incongruence phylogénétique
  - ▶ *Phylogénomique*
2. Artéfacts de reconstruction
  - ▶ *Causes et remèdes*
3. Transfert de gènes (EGT/HGT)
  - ▶ *Orthologie, paralogie et xénologie*
4. Prédiction phylogénétiques
  - ▶ *Horloges moléculaires*


# Les 6 Supergroupes

fusion des gènes DHFR et TS



'unikontes'  
vs. 'bikontes'





*Incongruence phylogénétique*  
*Phylogénomique*

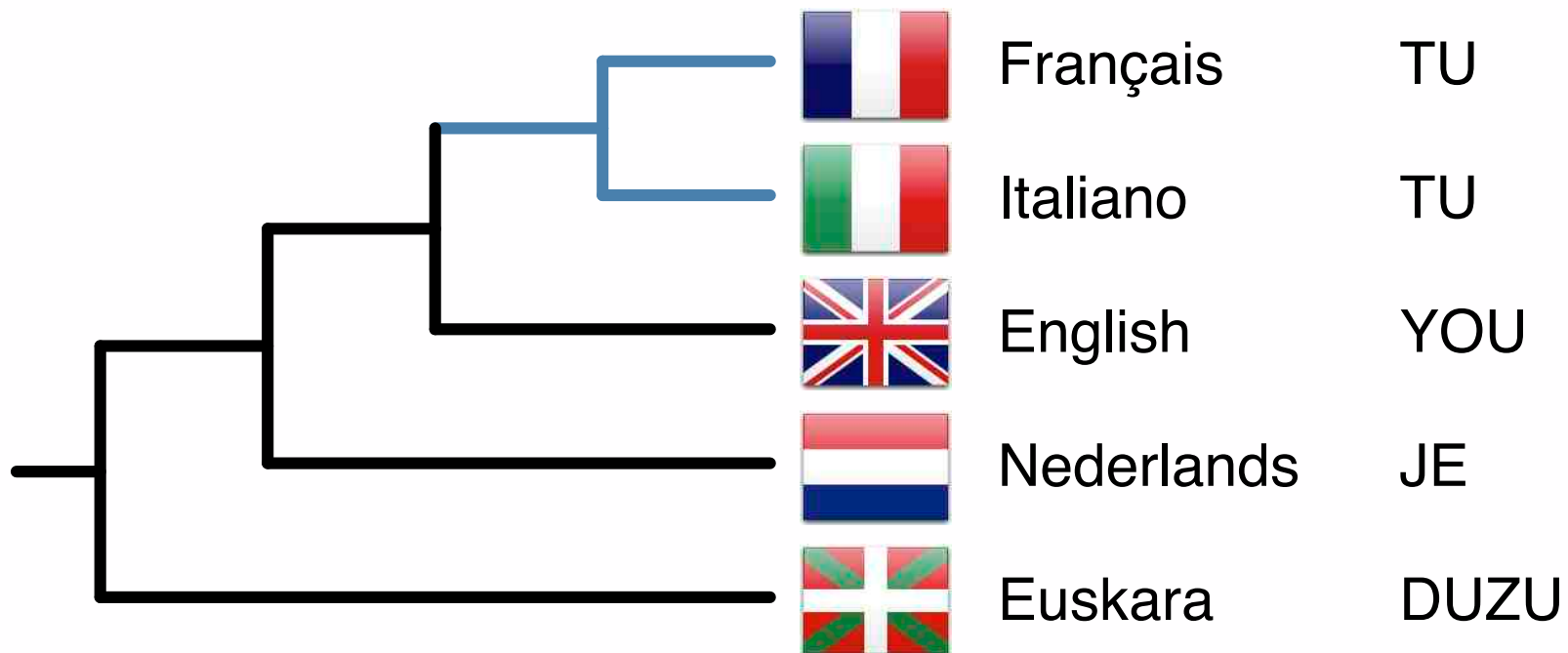


# Incongruence

un exemple tiré de la linguistique

## Potential relationships among European languages

based on the naive analysis of 1 word

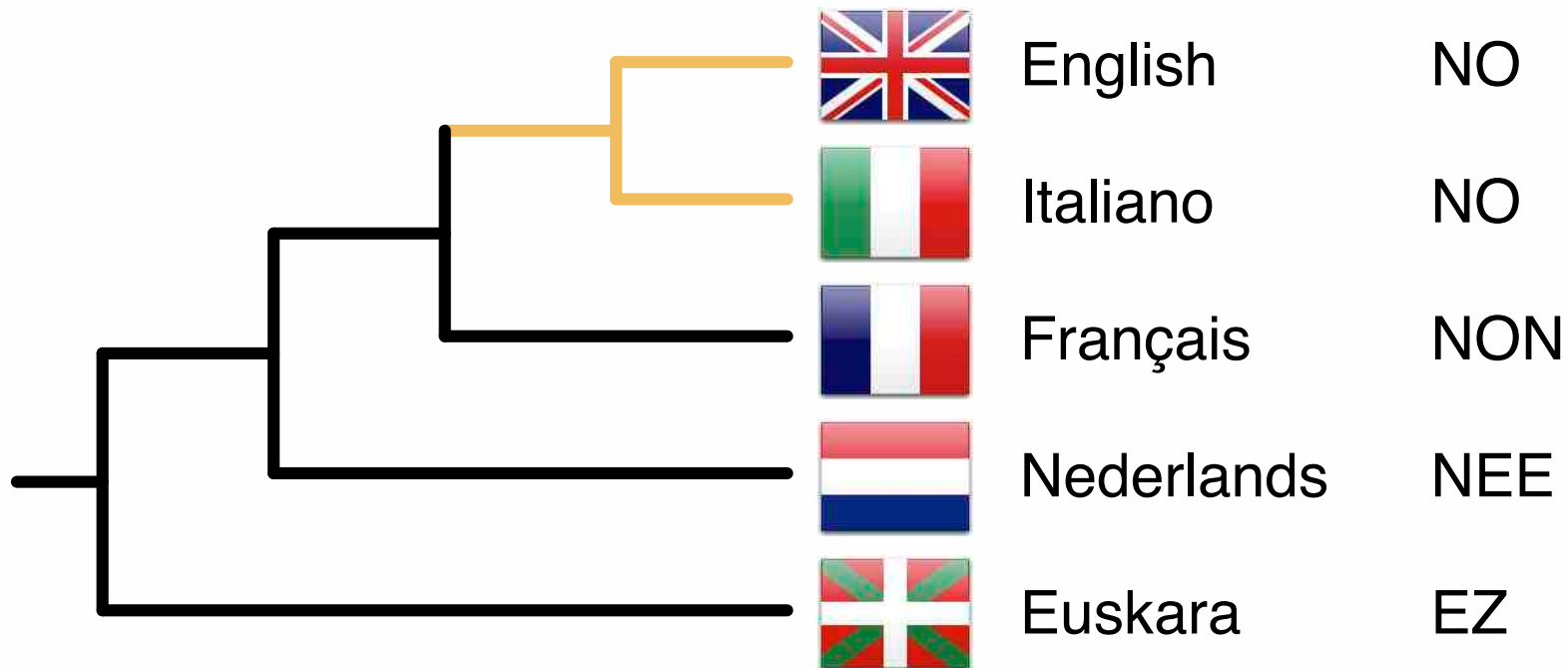


# Incongruence

un exemple tiré de la linguistique

## Potential relationships among European languages

based on the naive analysis of 1 word



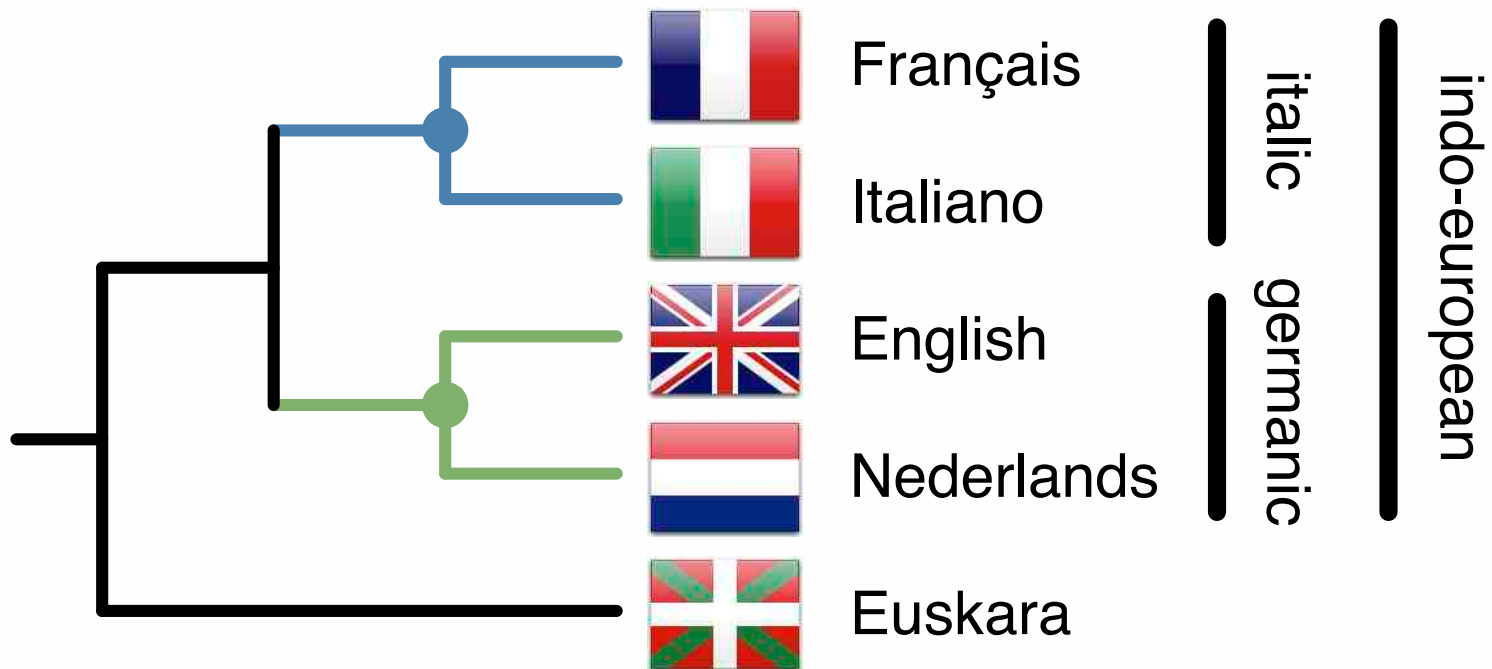
	français	italiano	english	nederlands	euskara
1	un	uno	one	een	bat
2	deux	due	two	twee	bi
3	trois	tre	three	drie	hiru
4	je	io	I	ik	ni
5	tu	tu	you	je	duzu
6	qui ?	chi?	who?	wie?	nor?
7	oui	si	yes	ja	bai
8	non	no	no	nee	ez
9	mère	madre	mother	moeder	ama
10	père	padre	father	vader	aita
11	dent	dente	tooth	tand	hortz
12	coeur	cuore	heart	hart	bihotza
13	pied	piede	foot	voet	oinez
14	souris	topolino	mouse	muis	saguaren

# Incongruence

un exemple tiré de la linguistique

## Known relationships among European languages

strongly supported by the naive analysis of 14 words



# Phylogénomique

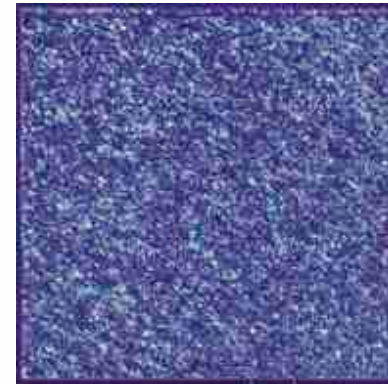
*La révolution génomique a commencé en 1995.*

## Genomics



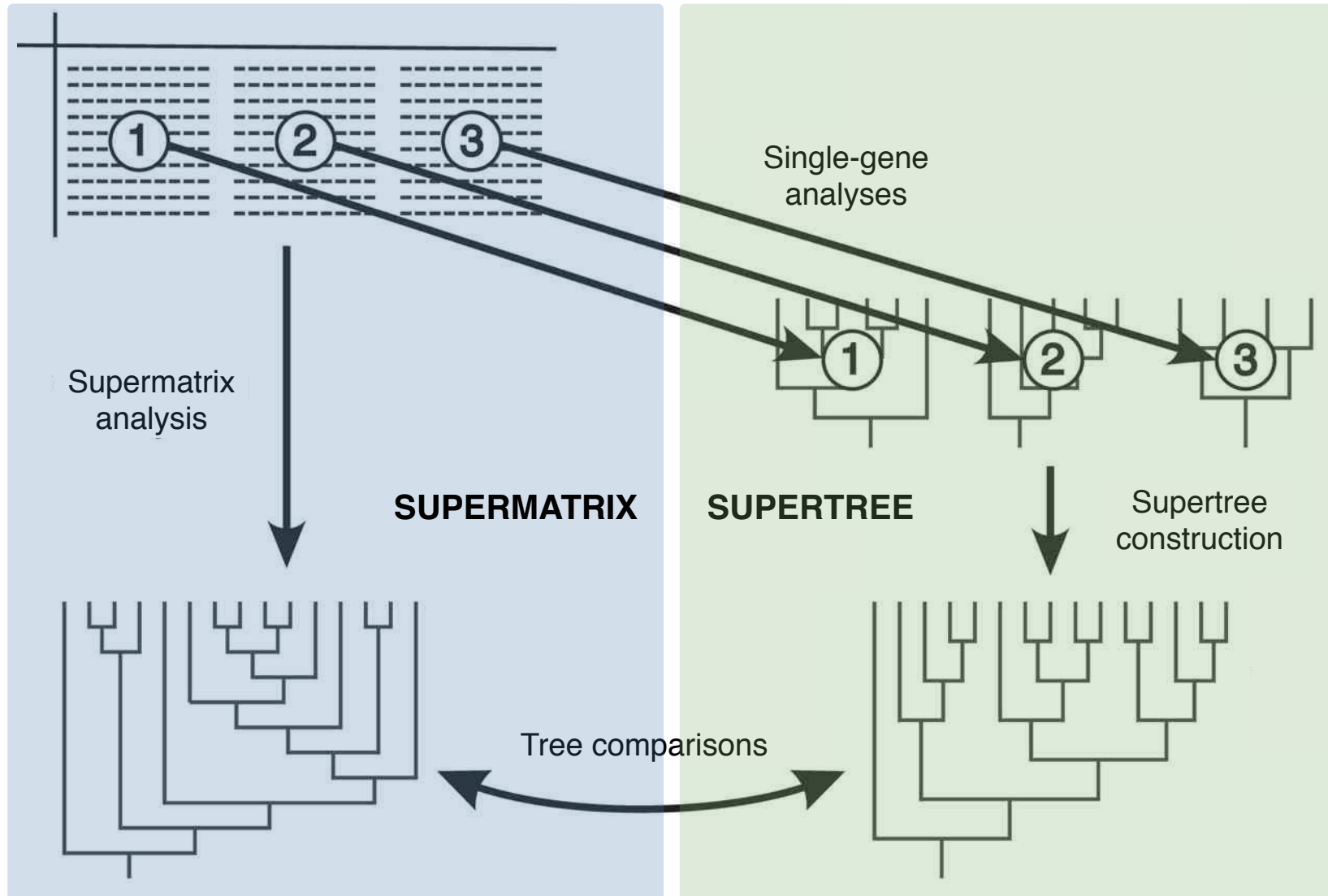
Automated Sanger  
sequencing machine  
(ABI 3730xl)

## High-throughput experiments



High-density  
oligonucleotide microarray  
(Affymetrix array)

# Phylogénomique



# Phylogénomique

## concaténation d'alignements en supermatrices

	$S^1$		$S^2$		$S^3$
$S_1$	A C G T C A A G		$S_1$ T G G - - T		$S_1$ C G G A C T A C G T
$S_2$	A C - T C C A G		$S_3$ A G C T C C		$S_4$ C C C T - - - - G G
$S_3$	A C - T C G A C		$S_4$ A G C T C G		$S_5$ C G T T C G A C G T

	$S^1$		$S^2$		$S^3$
$S_1$	A C G T C A A G		T G G - - T		C G G A C T A C G T
$S_2$	A C - T C C A G		. . . . .		. . . . .
$S_3$	A C - T C G A C		A G C T C C		. . . . .
$S_4$	. . . . .		A G C T C G		C C C T - - - - G G
$S_5$	. . . . .		. . . . .		C G T T C G A C G T

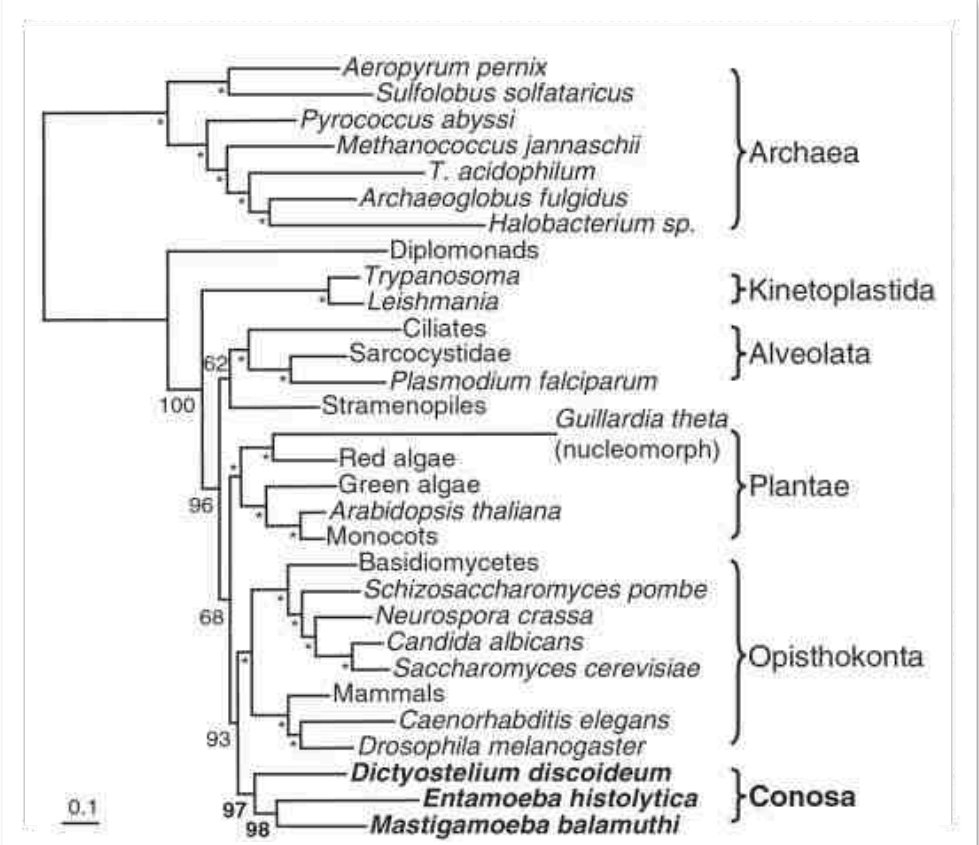
# The analysis of 100 genes supports the grouping of three highly divergent amoebae: *Dictyostelium*, *Entamoeba*, and *Mastigamoeba*

Eric Bapteste\*, Henner Brinkmann†, Jennifer A. Lee‡, Dorothy V. Moore‡, Christoph W. Sensen§, Paul Gordon¶, Laure Duruflé\*, Terry Gaasterland‡, Philippe Lopez\*, Miklós Müller‡, and Hervé Philippe\*||

\*Unité Mixte de Recherche 7622 Centre National de la Recherche Scientifique, Université Paris 6, 9 Quai Saint Bernard, Bât C, 75005 Paris, France; †Department of Biology, University of Konstanz, 78457 Konstanz, Germany; ‡The Rockefeller University, 1230 York Avenue, New York, NY 10021; §Department of Biochemistry and Molecular Biology, University of Calgary, 3330 Hospital Drive N.W., Calgary, AB, Canada, T2N 4N1; and ¶National Research Council, Institute for Marine Biosciences, 1411 Oxford Street, Halifax, NS, Canada B3H 3Z1

Communicated by William Trager, The Rockefeller University, New York, NY, December 11, 2001 (received for review October 10, 2001)

The phylogenetic relationships of amoebae are poorly resolved. To address this difficult question, we have sequenced 1,280 expressed sequence tags from *Mastigamoeba balamuthi* and assembled a large data set containing 123 genes for representatives of three phenotypically highly divergent major amoeboid lineages: Pelobionta, Entamoebidae, and Mycetozoa. Phylogenetic reconstruction was performed on  $\approx 25,000$  aa positions for 30 species by using maximum-likelihood approaches. All well-established eukaryotic groups were recovered with high statistical support, validating our approach. Interestingly, the three amoeboid lineages strongly clustered together in agreement with the Conosa hypothesis [as defined by T. Cavalier-Smith (1998) *Biol. Rev. Cambridge Philos. Soc.* 73, 203–266]. Two amitochondriate amoebae, the free-living *Mastigamoeba* and the human parasite *Entamoeba*, formed a significant sister group to the exclusion of the mycetozoan *Dictyostelium*. This result suggested that a part of the reductive process in the evolution of *Entamoeba* (e.g., loss of typical mitochondria) occurred in its free-living ancestors. Applying this inexpensive expressed sequence tag approach to many other lineages will surely improve our understanding of eukaryotic evolution.

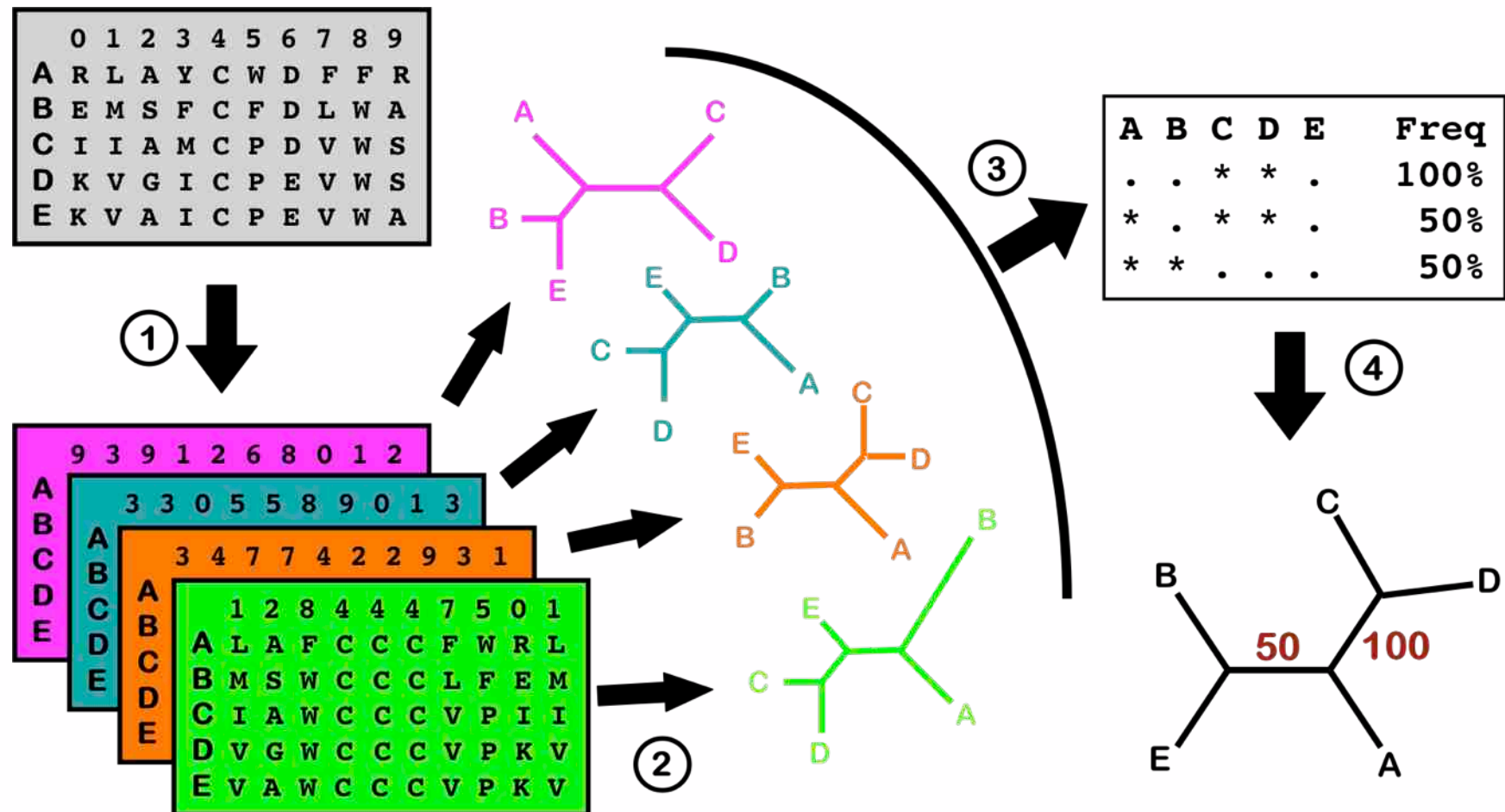


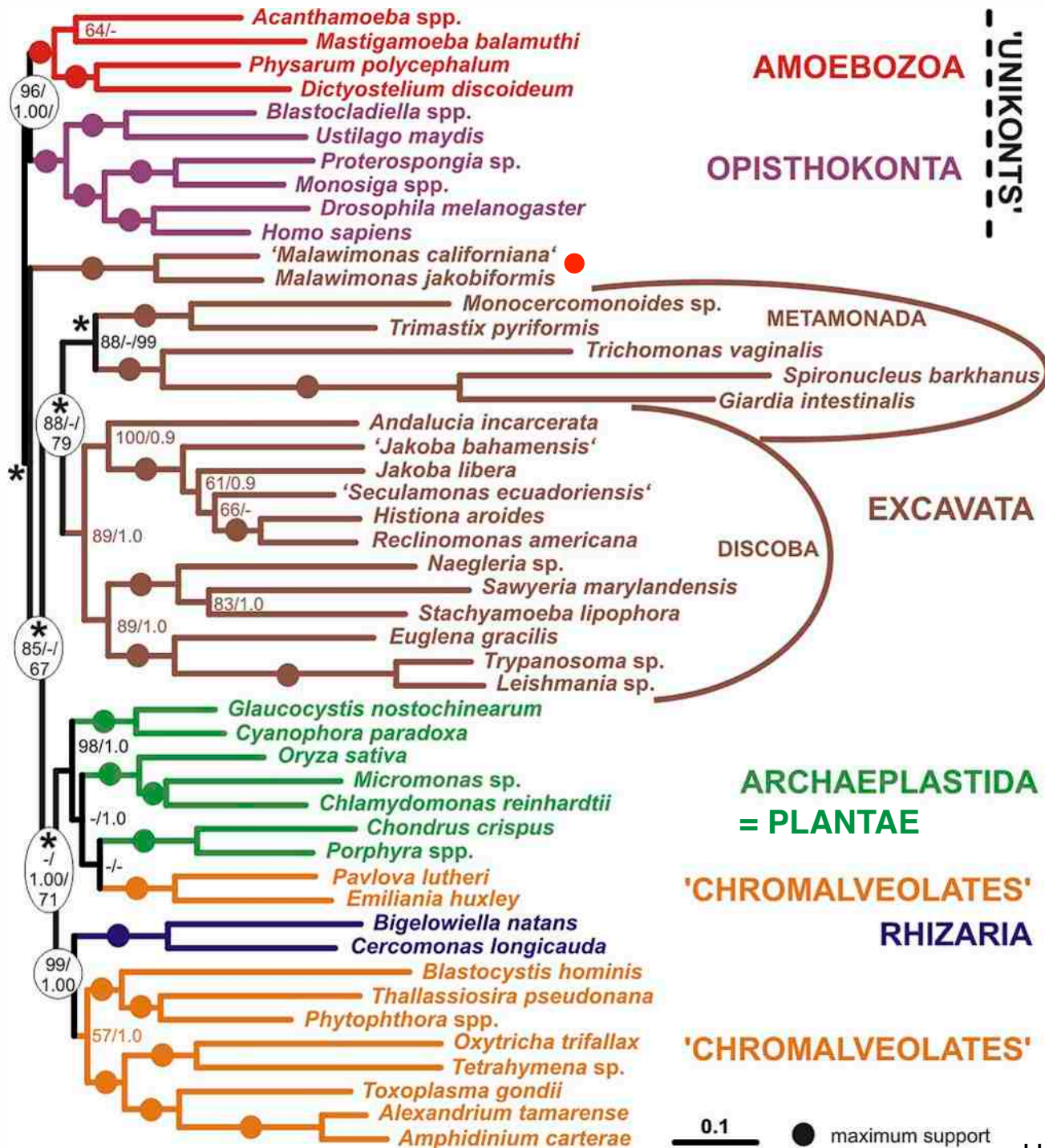
unicellular amoebae are possibly the simplest eukaryotic



# Phylogénomique

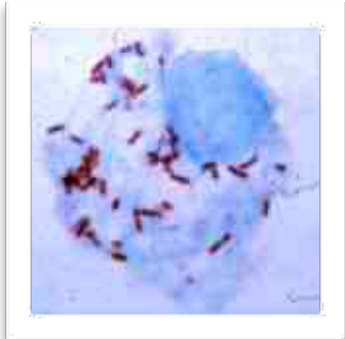
principe du bootstrap et robustesse des arbres





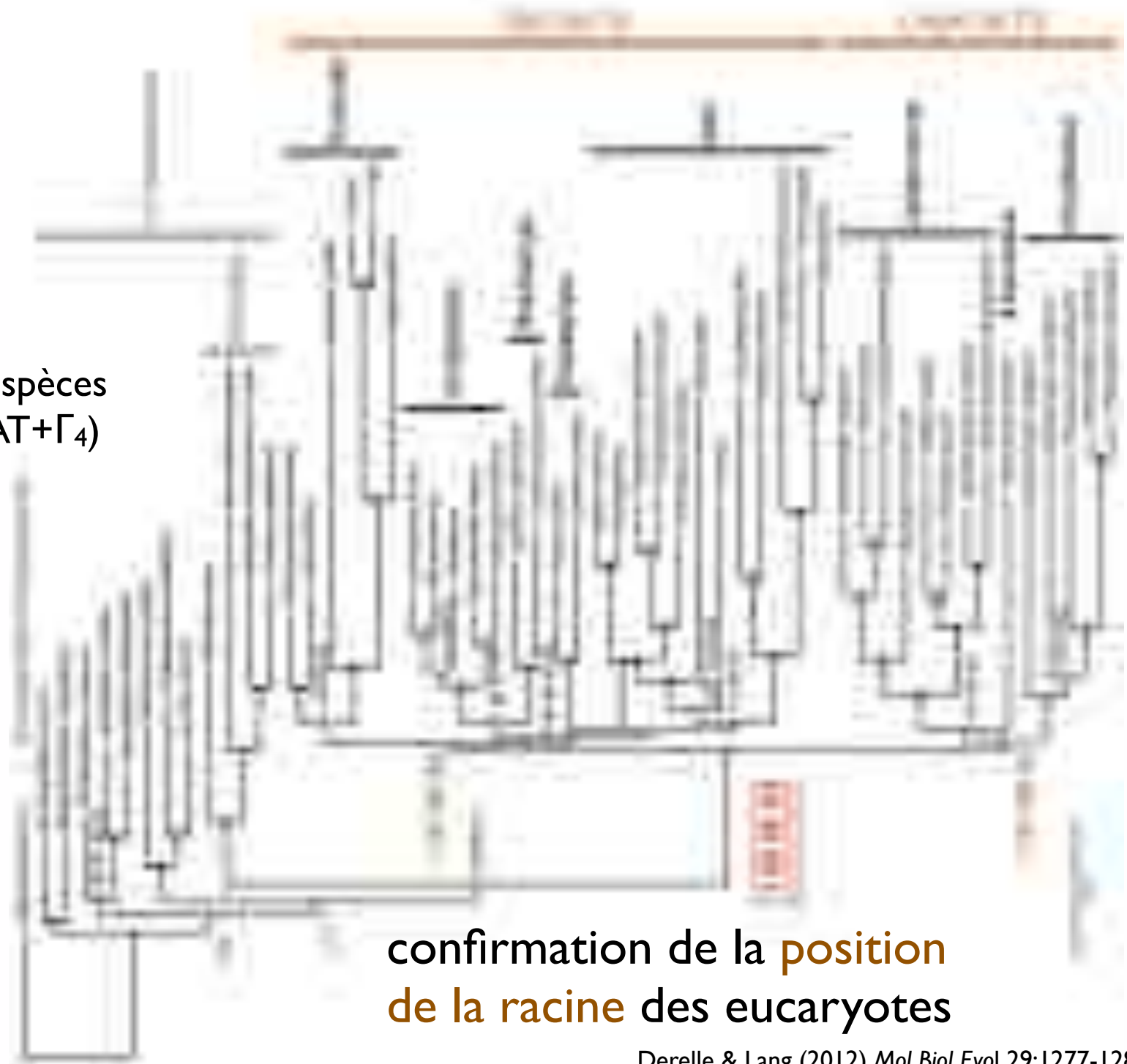
confirmation  
de l'existence  
de certains des  
supergroupes

143 gènes x 48 espèces  
RAXML (WAG+Γ<sub>4</sub>)



*Rickettsia*

42 gènes x 54 espèces  
PhyloBayes (CAT+ $\Gamma_4$ )



confirmation de la **position**  
de la **racine** des eucaryotes



*Artéfacts de reconstruction*  
*Causes et remèdes*

# Artéfacts

une question de substitutions multiples

# Artéfacts

une question de substitutions multiples



# Artéfacts

une question de substitutions multiples

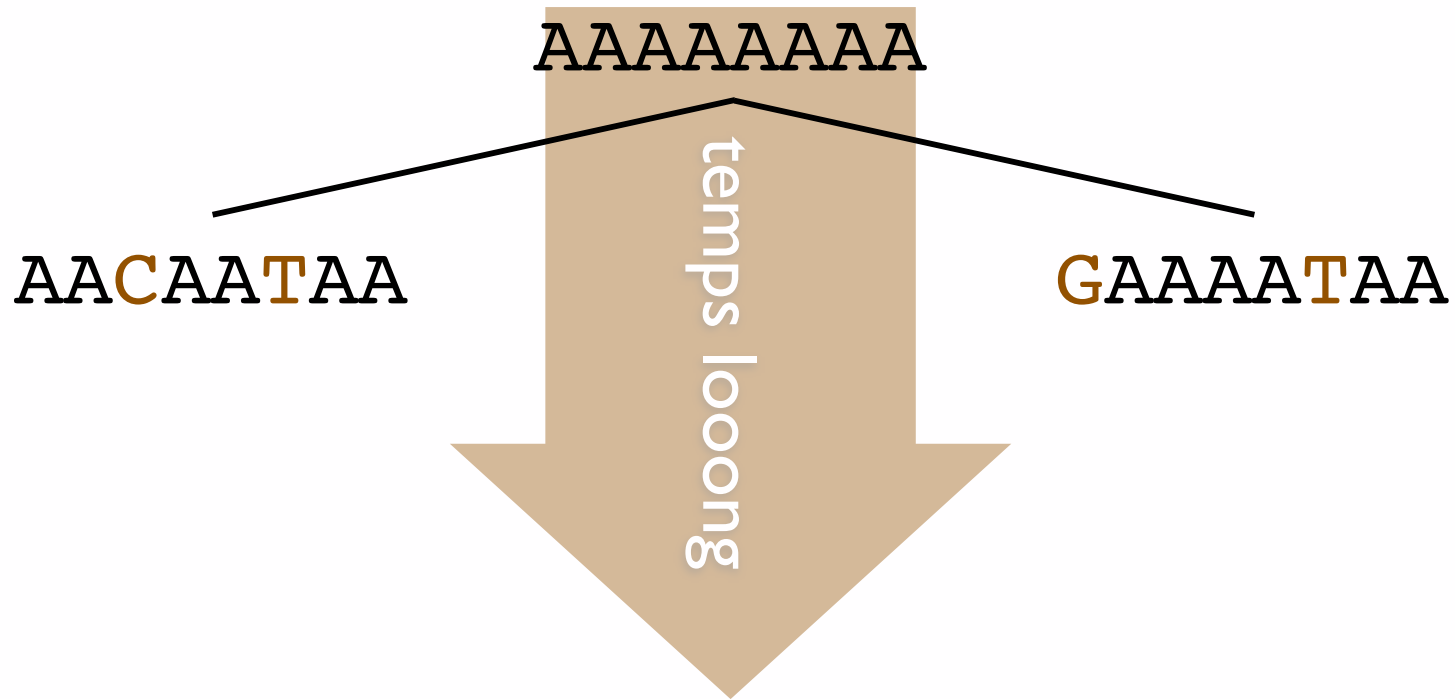
AAAAAAAAAA



temps looong

# Artéfacts

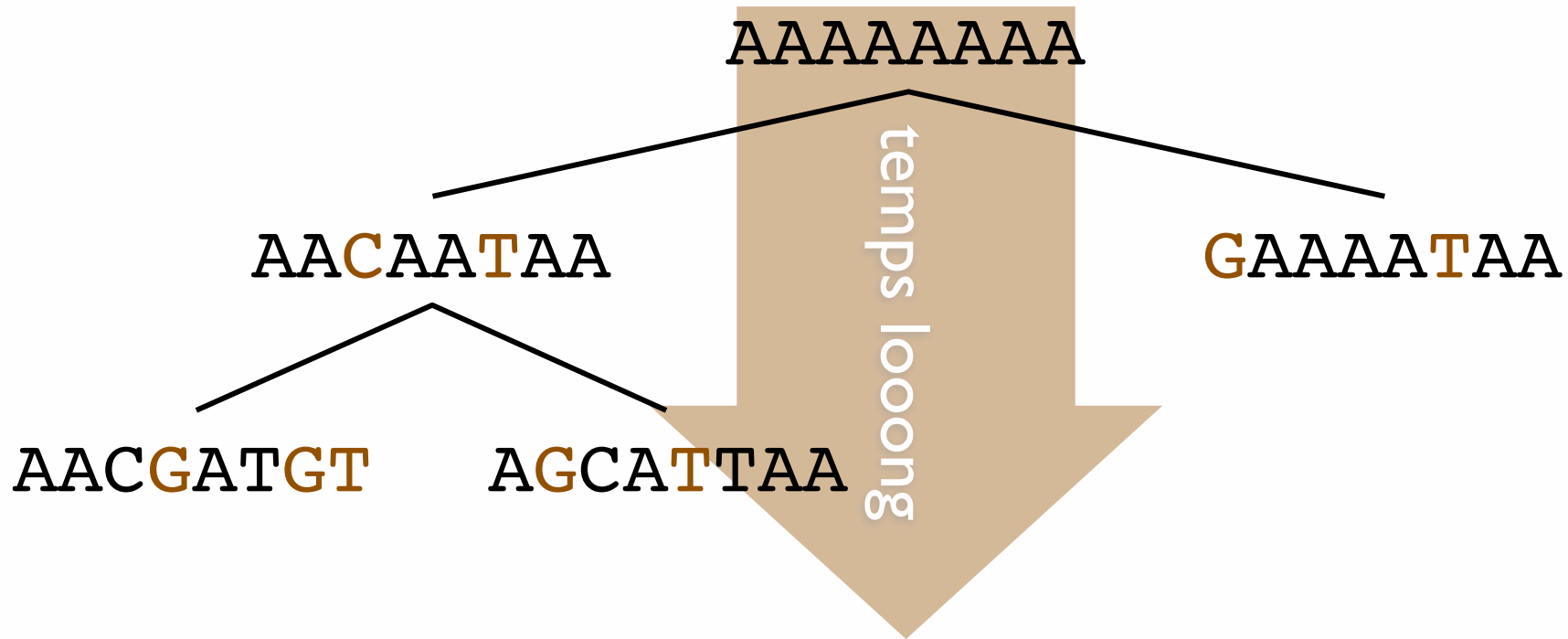
une question de substitutions multiples





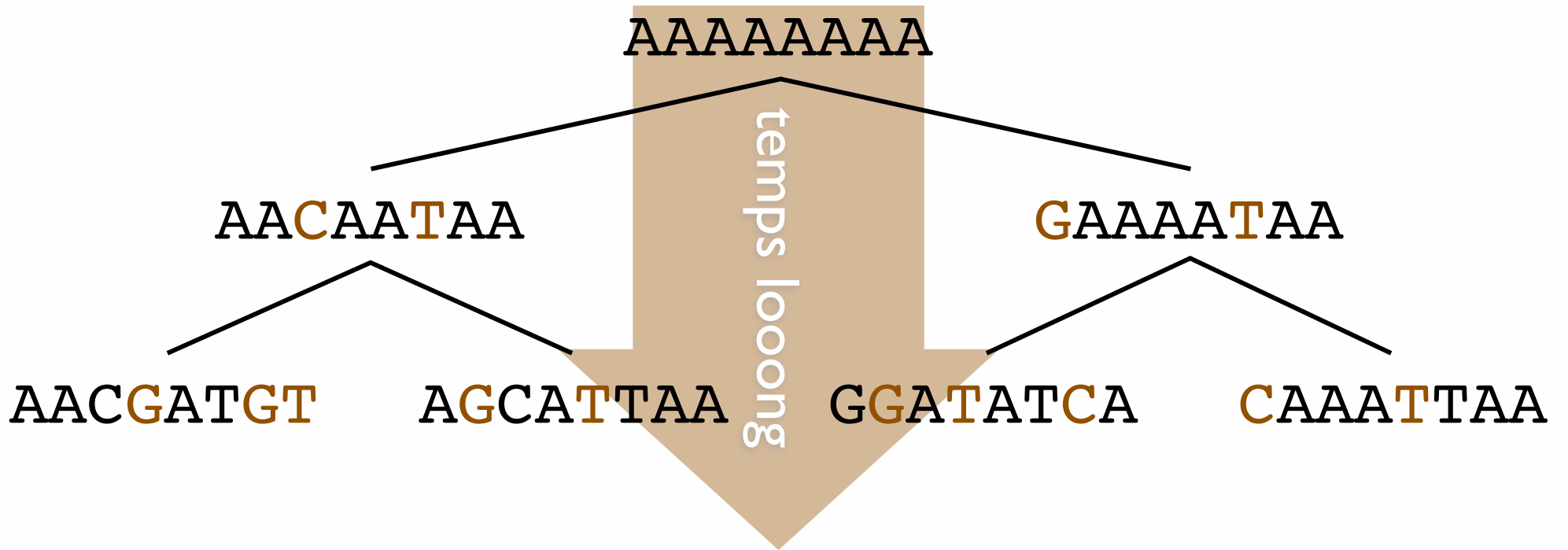
# Artéfacts

une question de substitutions multiples



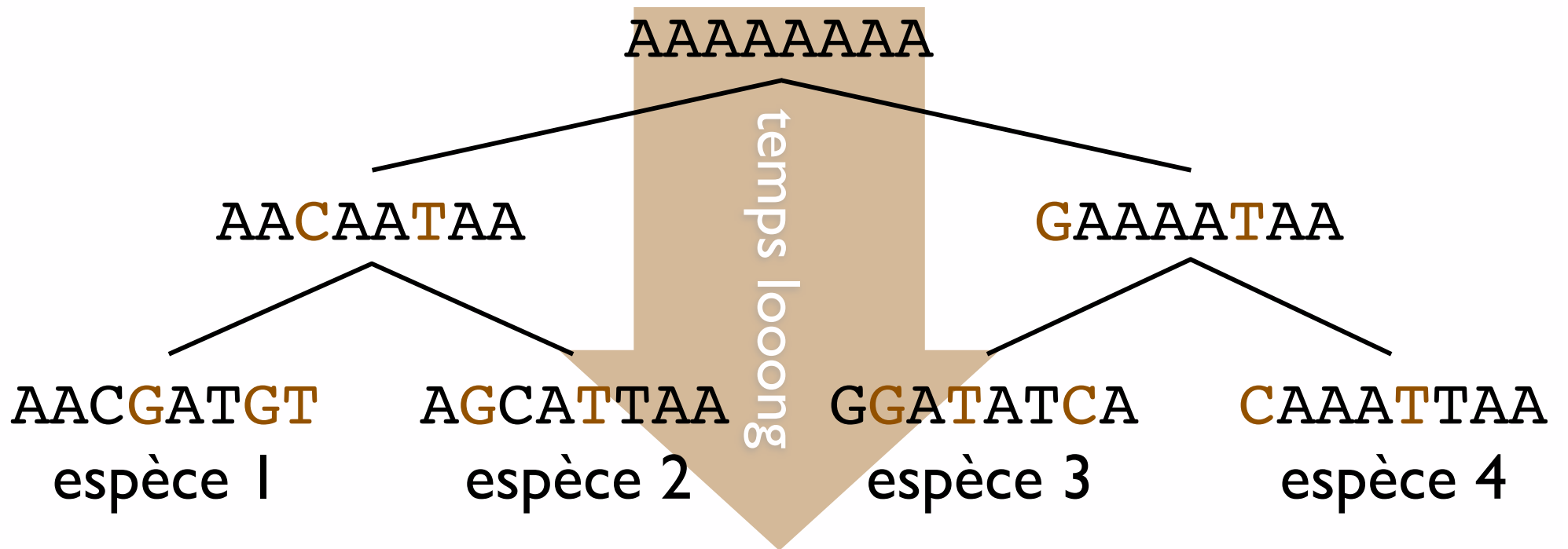
# Artéfacts

une question de substitutions multiples



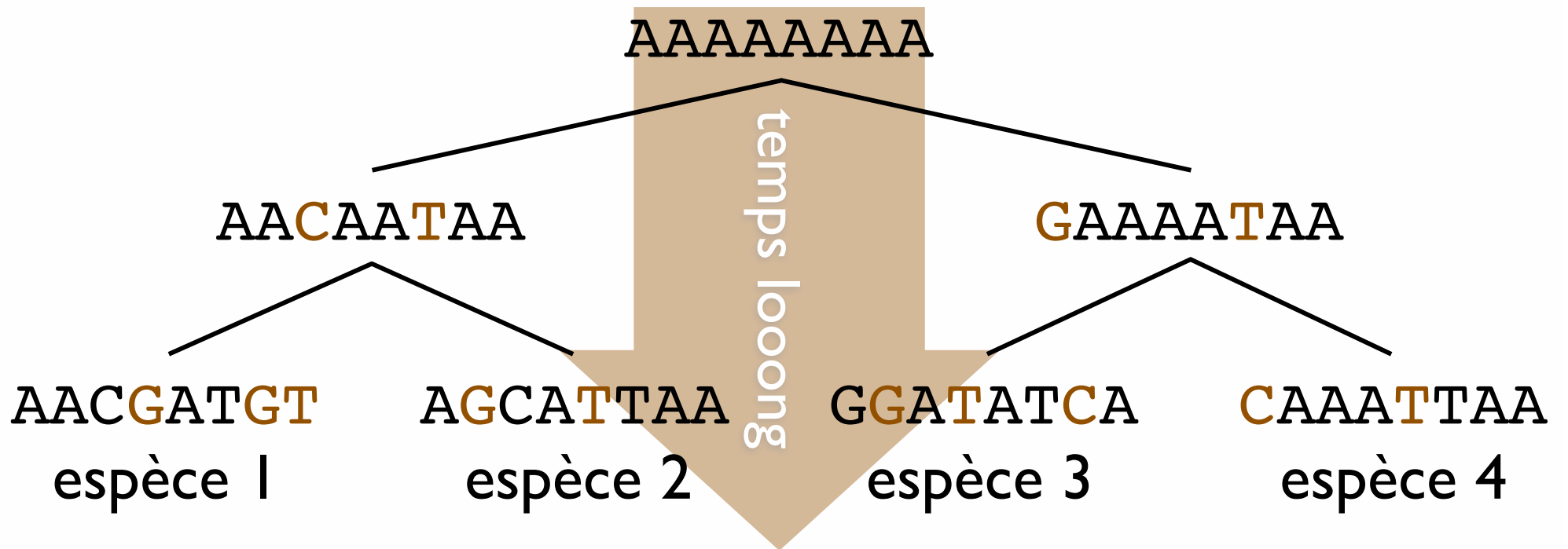
# Artéfacts

une question de substitutions multiples



# Artéfacts

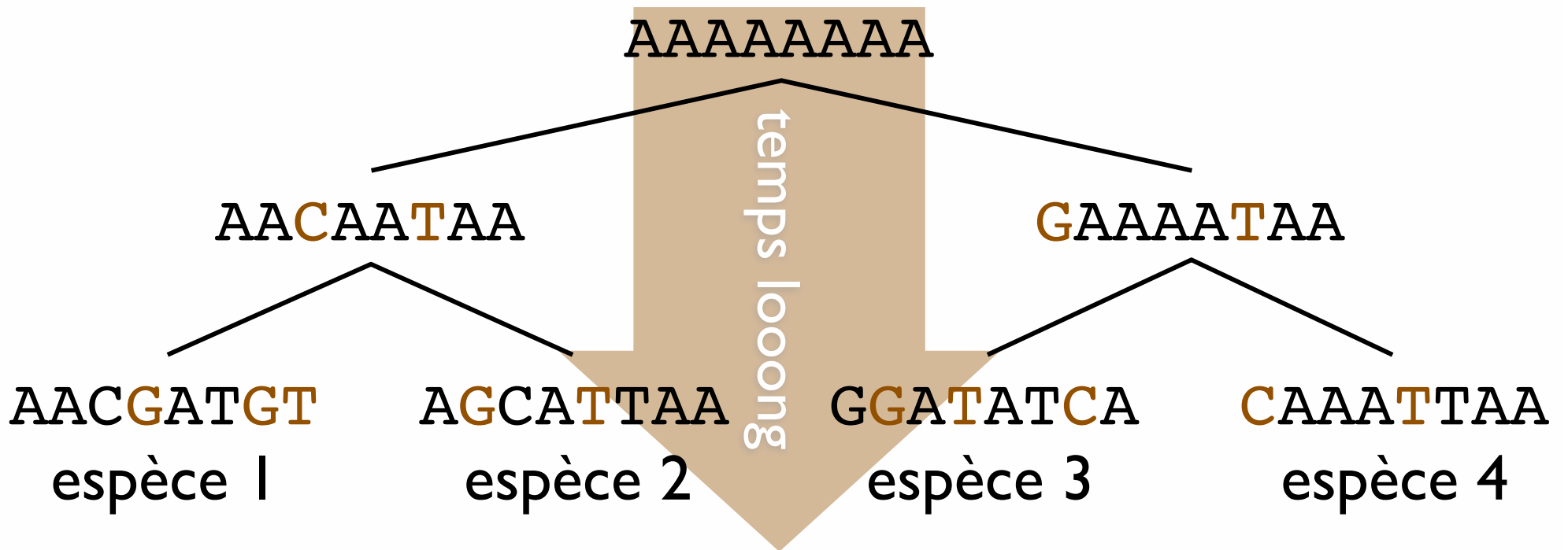
une question de substitutions multiples



- 1 AACGATGT
- 2 AGCATTAA
- 3 GGATATCA
- 4 CAAATTAA

# Artéfacts

une question de substitutions multiples



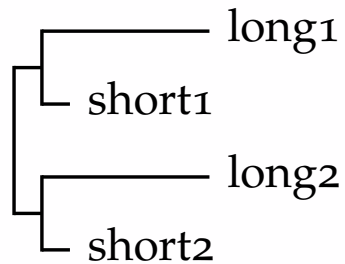
- 1 AACGATGT
- 2 AGCATTAA
- 3 GGATATCA
- 4 CAAATTAA

Lorsque l'homoplasie est trop importante, les sites informatifs deviennent trompeurs.

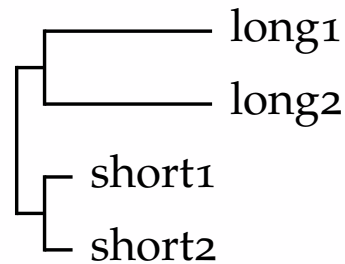
C'est le signal non-phylogénétique.

# Artéfacts

## attraction des longues branches / inconsistance

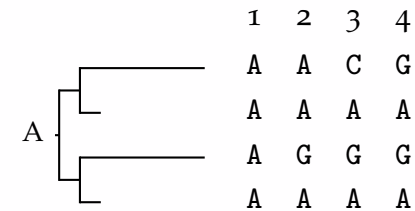


True



Recovered

- 1 both the same as short1 and short2
- 2 one the same and one different
- 3 both different and different from each other
- 4 both different but the same as each other



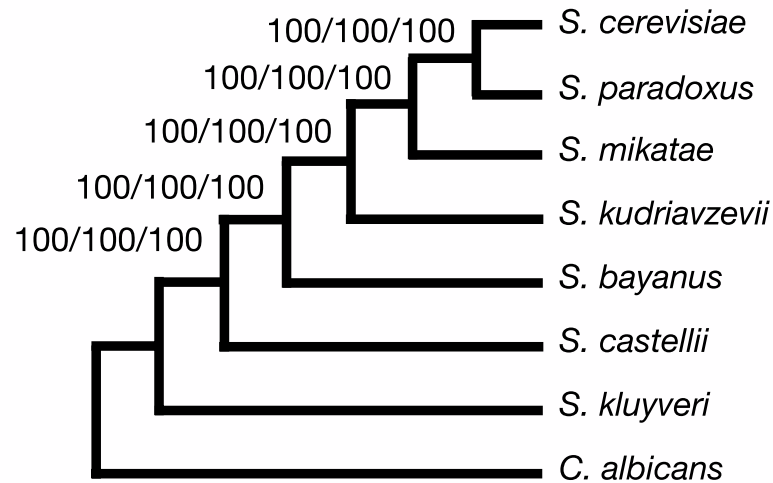
- |   |      |                          |
|---|------|--------------------------|
| 1 | AAAA | parsimony uninformative  |
| 2 | AAGA | parsimony uninformative  |
| 3 | CAGA | parsimony uninformative  |
| 4 | GAGA | parsimony misinformative |

Puisque tous les sites sont **soit non-informatifs, soit trompeurs**, la parcimonie tend à regrouper les longues branches ensemble, et ce d'autant plus solidement qu'elle dispose de beaucoup de données.

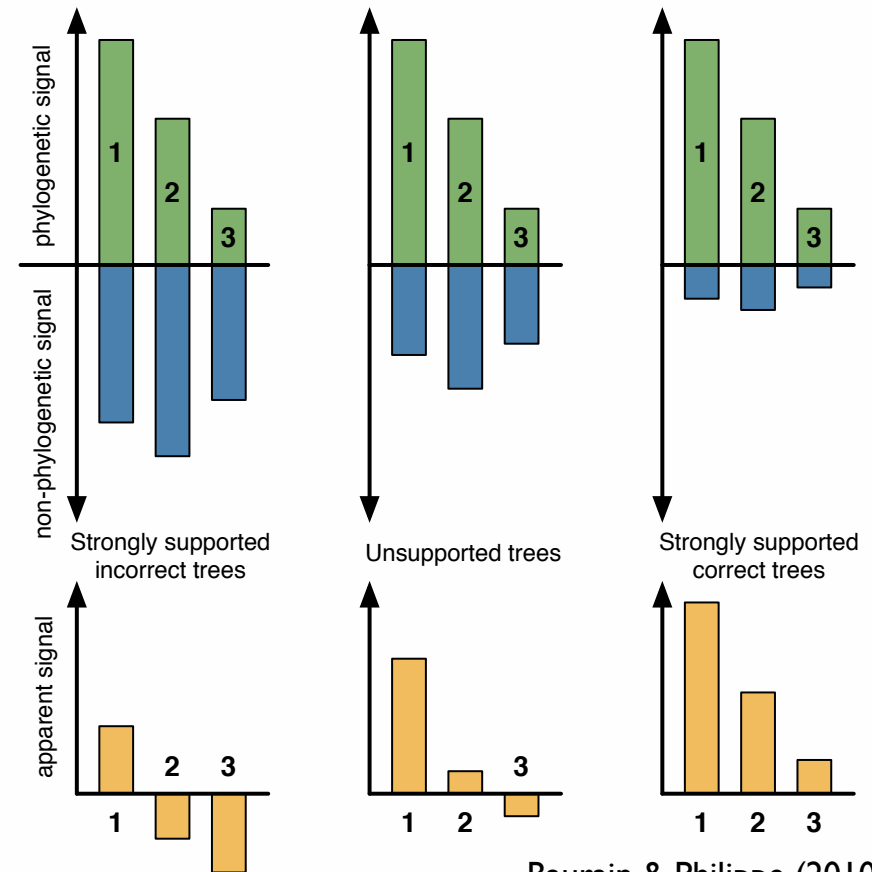
# Ending incongruence

Henry Gee

Recovering the true evolutionary history of any group of organisms has seemed impossible. The availability of large amounts of genomic data promises an era in which the uncertainties are better constrained.



Rokas et al. (2003) *Nature* 425:798-804



Baurain & Philippe (2010)  
In *Evolutionary Genomics and Systems Biology*



Jeffroy et al. (2006) *Trends Genet* 22:225-231



# Phylogenomics: the beginning of incongruence?

Olivier Jeffroy, Henner Brinkmann, Frédéric Delsuc and Hervé Philippe

Canadian Institute for Advanced Research, Centre Robert-Cedergren, Département de Biochimie, Université de Montréal, Succursale Centre-Ville, Montréal, Québec, Canada, H3C3J7

# Artéfacts

*Comment réduire le signal non-phylogénétique ?*

★ Combinaison de 3 approches :

1. Amélioration de l'échantillonnage taxonomique

*a. Choix d'espèces à évolution lente*

*b. Augmentation de la densité en espèces*

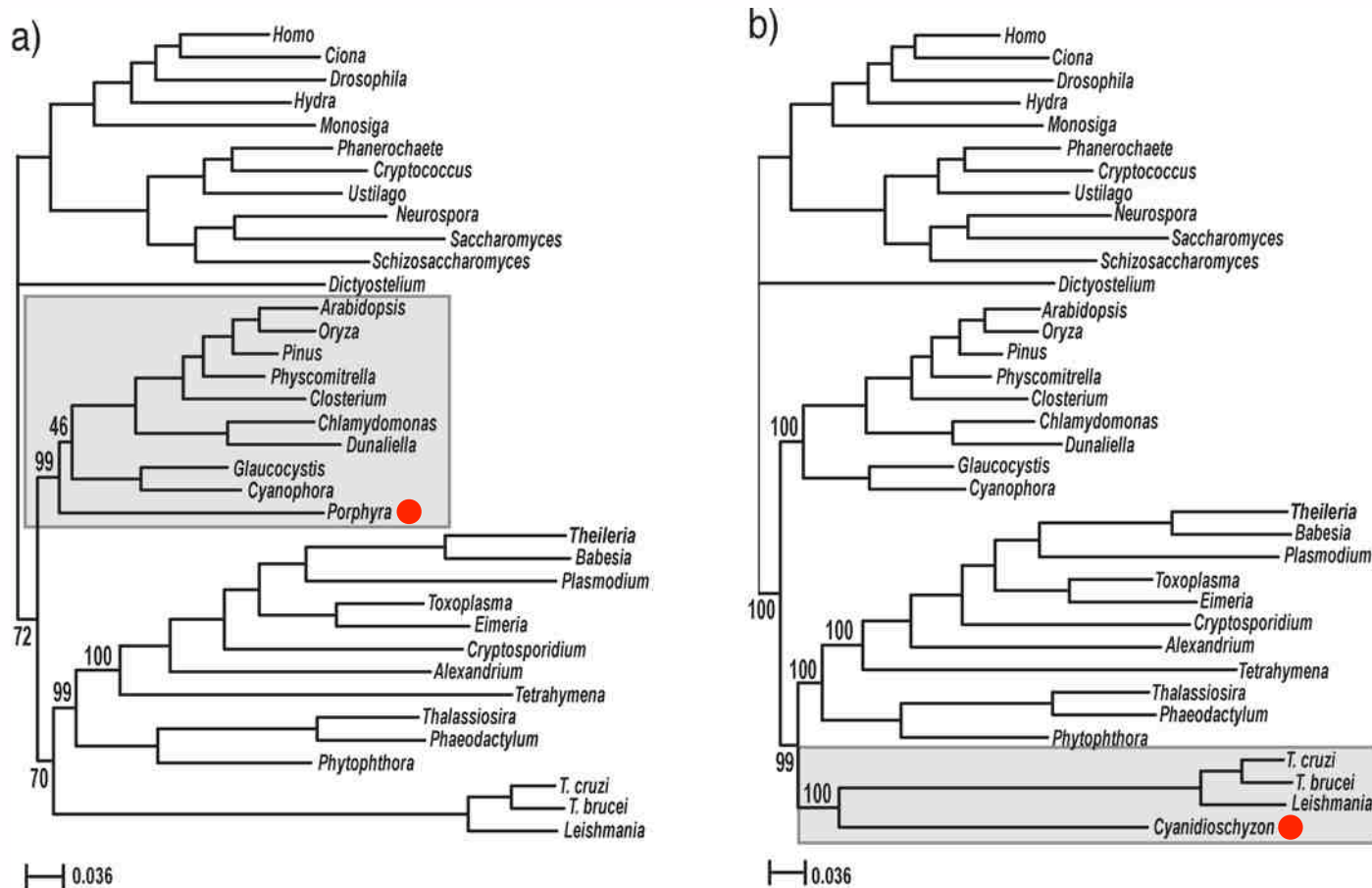
2. Elimination des sites à évolution rapide

3. Utilisation de modèles d'évolution sophistiqués



# Remèdes

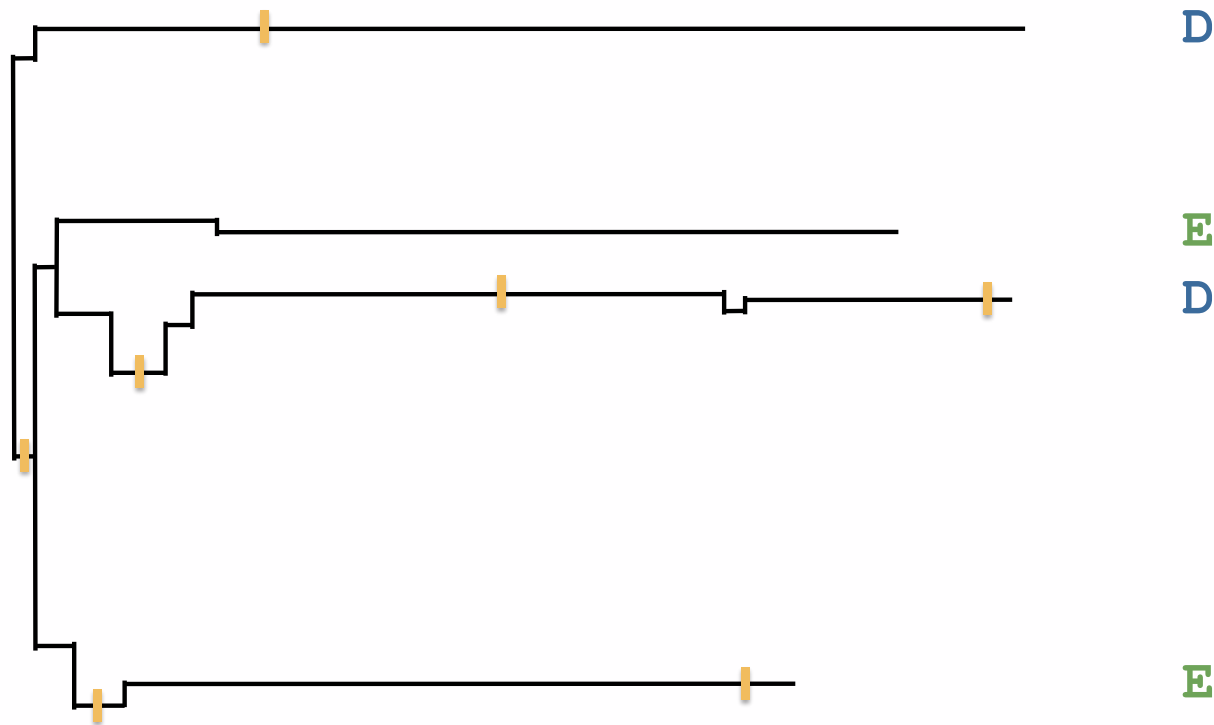
## Ia. choix d'espèces à évolution lente



Le bootstrap évalue l'homogénéité du signal et non sa nature.

# Remèdes

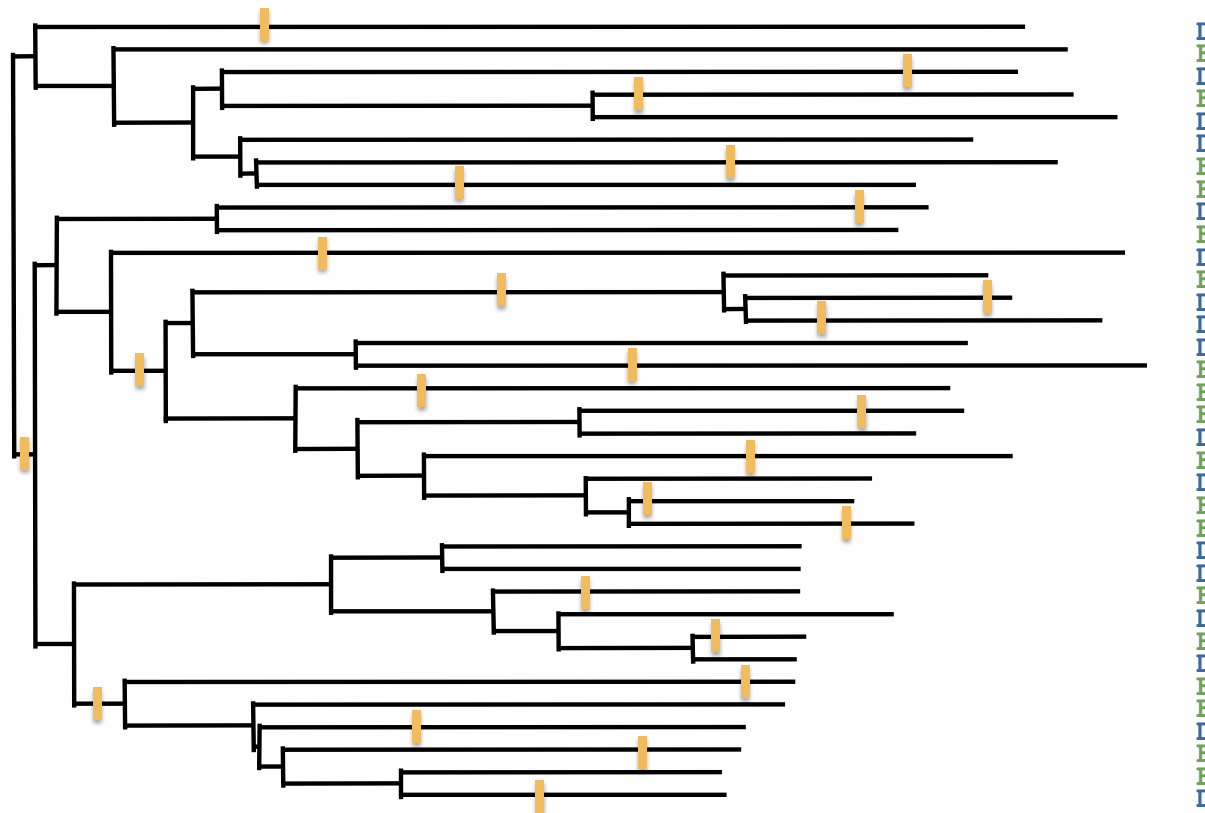
## I b. augmentation de la densité en espèces



Un échantillonnage épars masque les substitutions multiples.

# Remèdes

## I b. augmentation de la densité en espèces

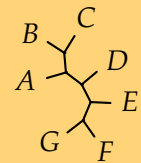


Un échantillonnage dense révèle les substitutions multiples.

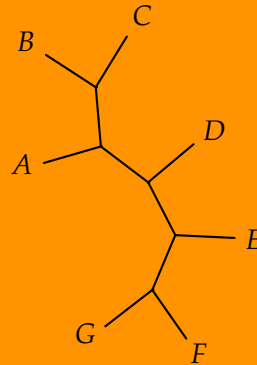
# Remèdes

## 2. élimination des sites à évolution rapide

A ttctaacgggacagtgcgcccactcacgcacctgggtcactgtatgcgagt  
B tgpgaaggtgctattgpgagcattcacgcagatggtaactgtatgtgaga  
C tgpgaaggtggtattgcccacattcgpgcggaaggtaacactatgtgaga  
D tcccatagpgacatggpgcatactgactcctatggatactgtatgpgagt  
E tgpgaaggpgacattgpgcacattcacgcataatggttactgtacgtgagt  
F tgpgaaggpggcattgpgtccattcacccgcatggagactttatgtgaga



Tree made using  
slow sites only

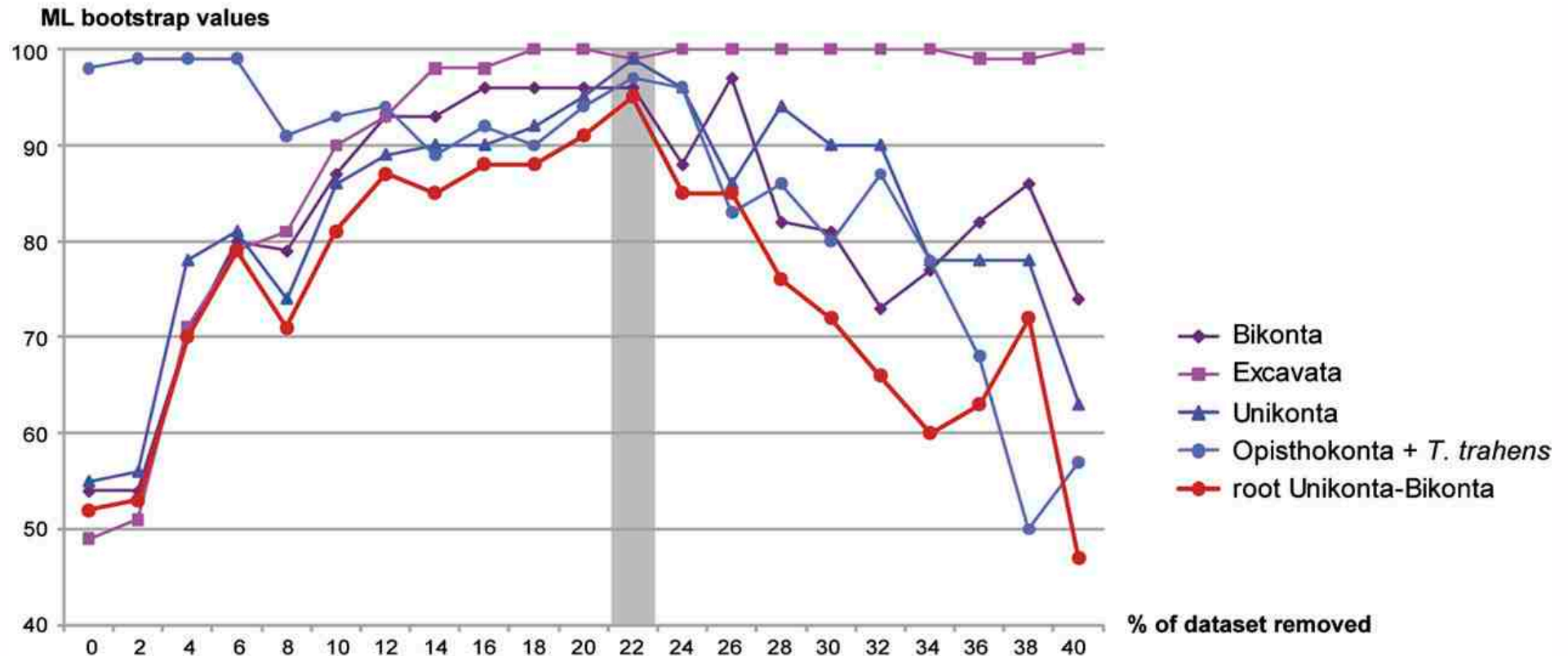


Tree made using  
fast sites only

Les sites sous faible  
sélection sont plus  
**saturés** que ceux sous  
forte sélection.

# Remèdes

## 2. élimination des sites à évolution rapide

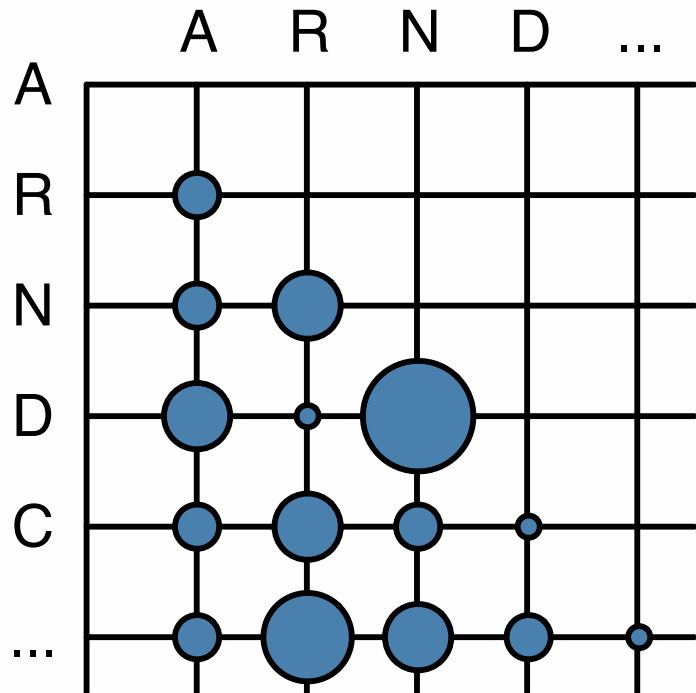


Le retrait des sites rapides **améliore la robustesse** de l'arbre.

# Remèdes

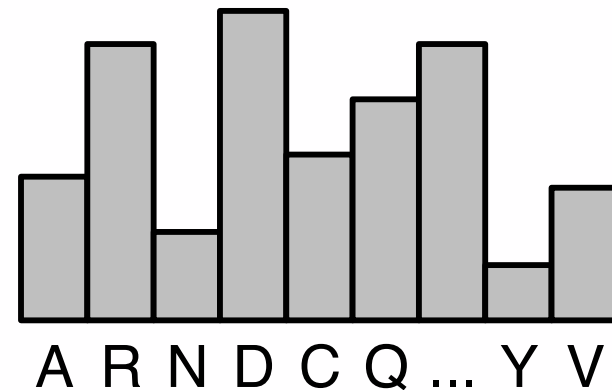
## 3. utilisation de modèles d'évolution sophistiqués

LG or WAG



1 global precomputed  
replacement matrix

+



1 global  
compositional profile

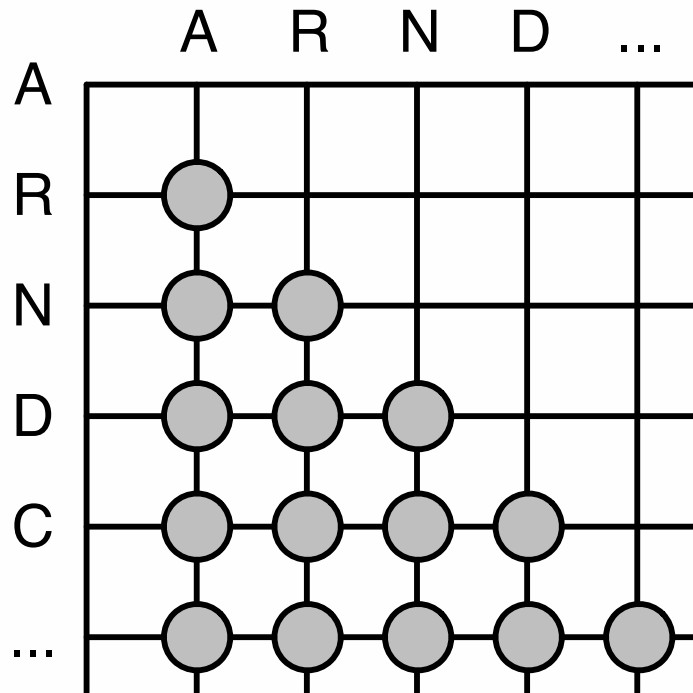
# Remèdes

## 3. utilisation de modèles d'évolution sophistiqués

CAT

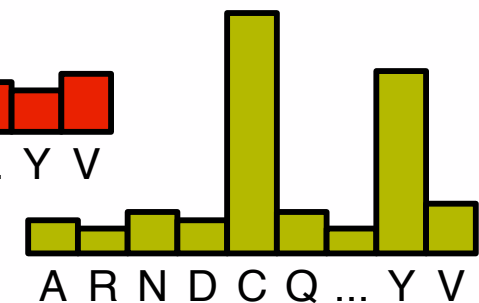
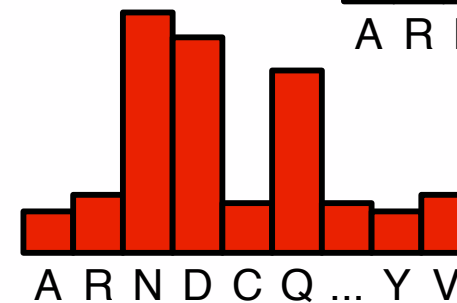
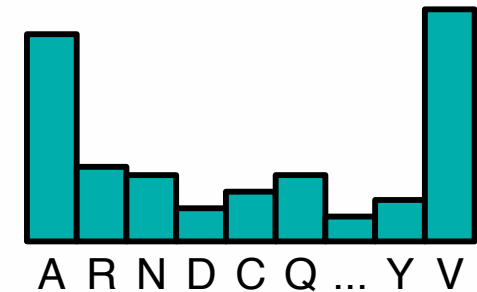


Nicolas Lartillot



1 global « flat »  
replacement matrix

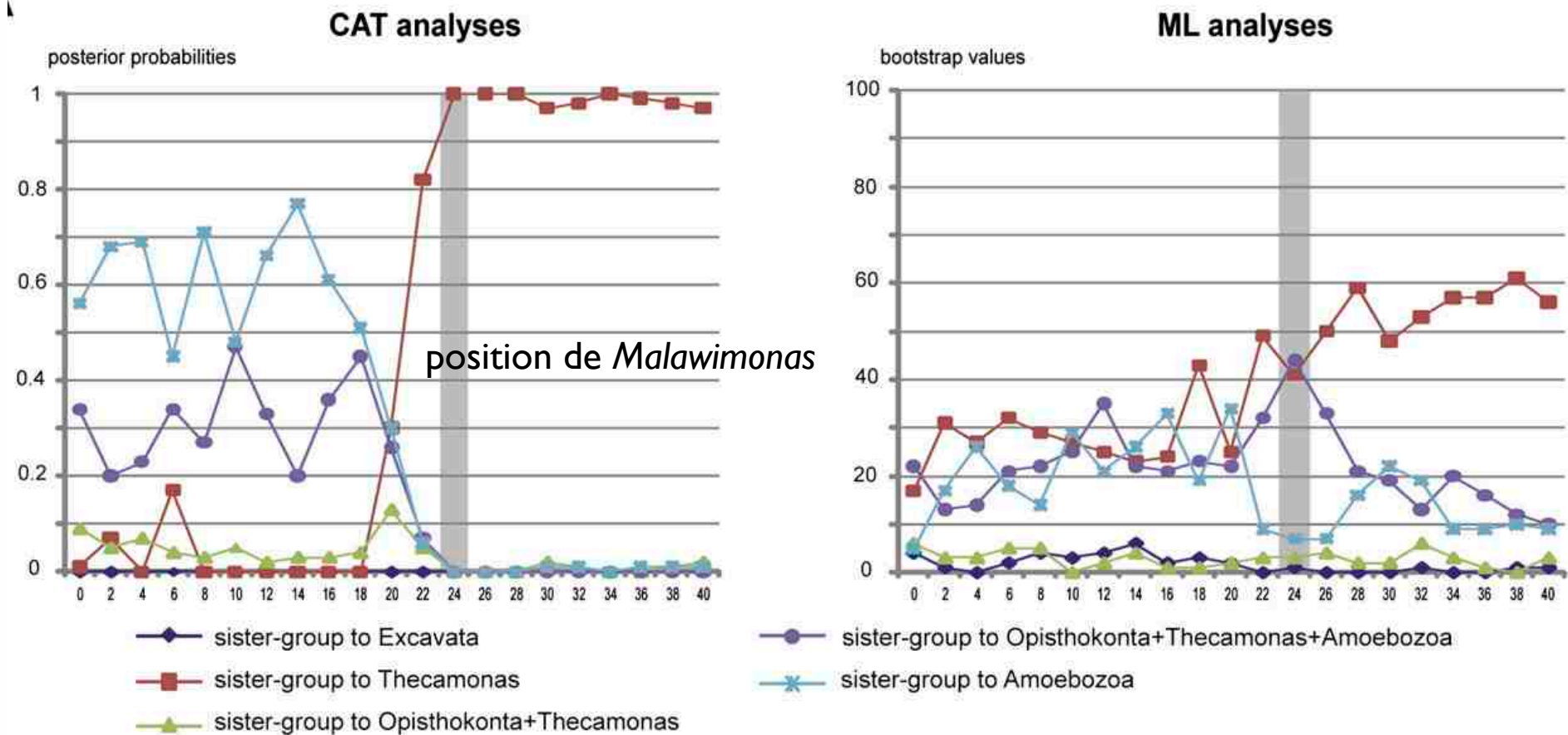
+



$K$  distinct  
compositional profiles

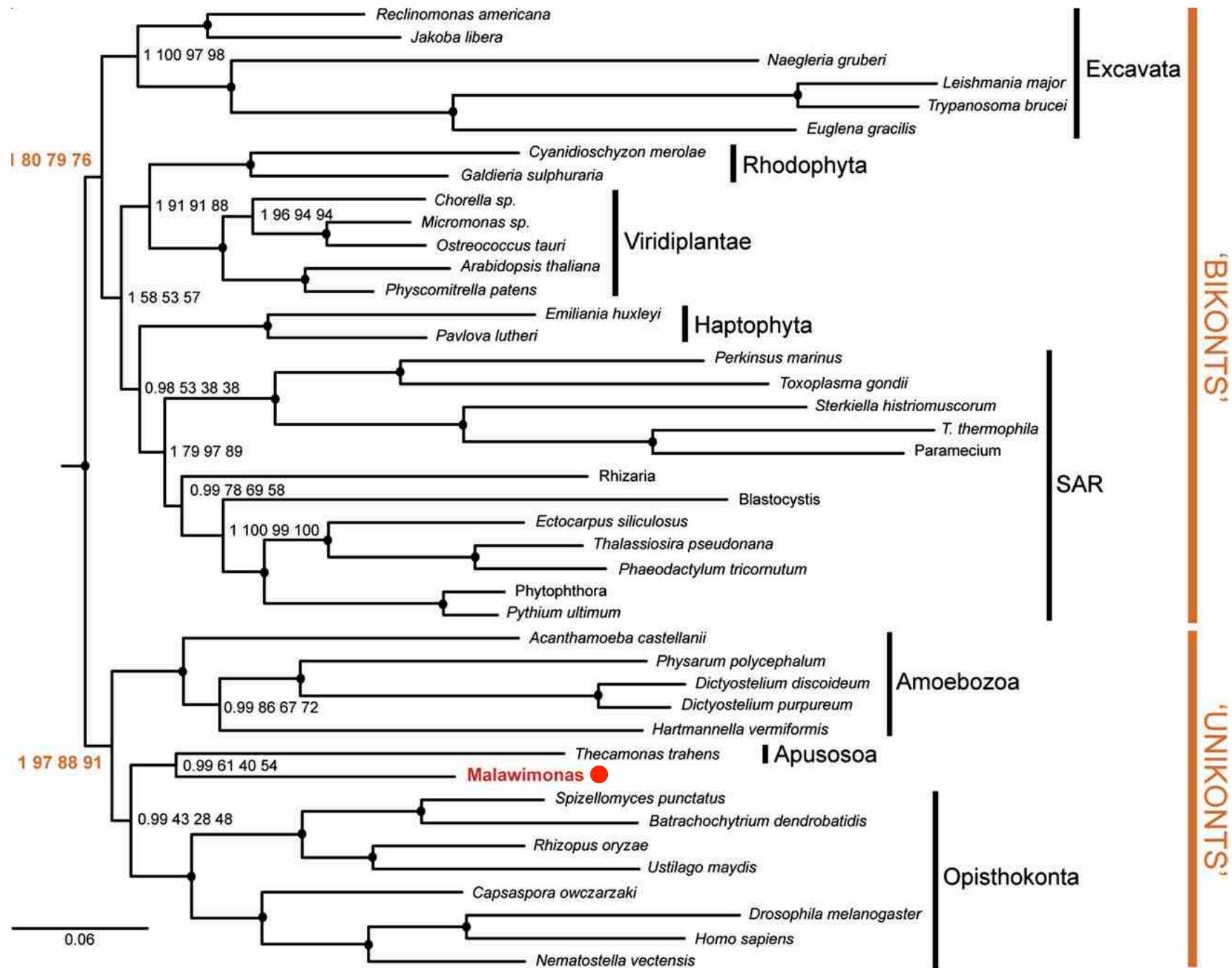
# Remèdes

## 3. utilisation de modèles d'évolution sophistiqués

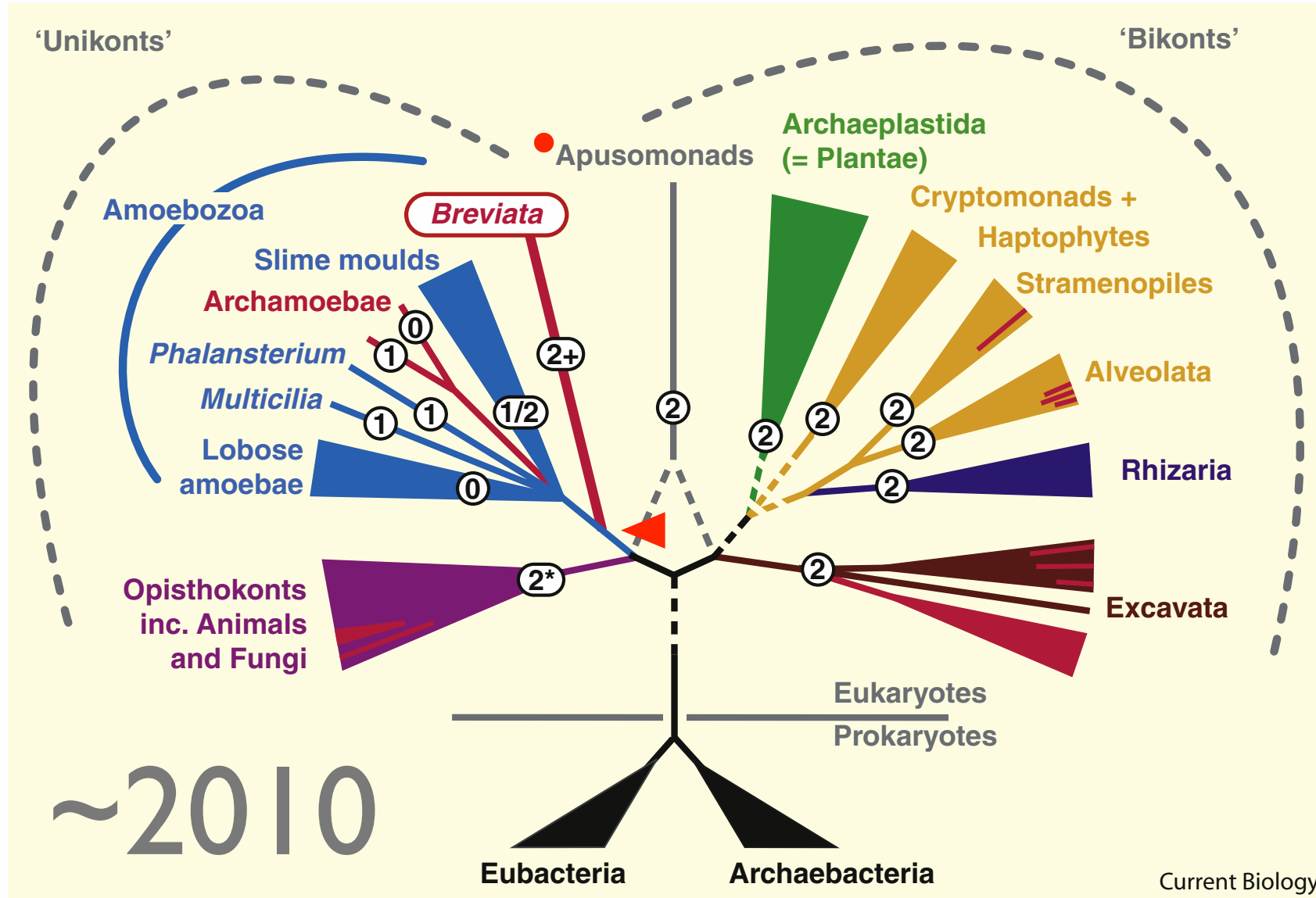


**CAT améliore la détection des substitutions multiples.**

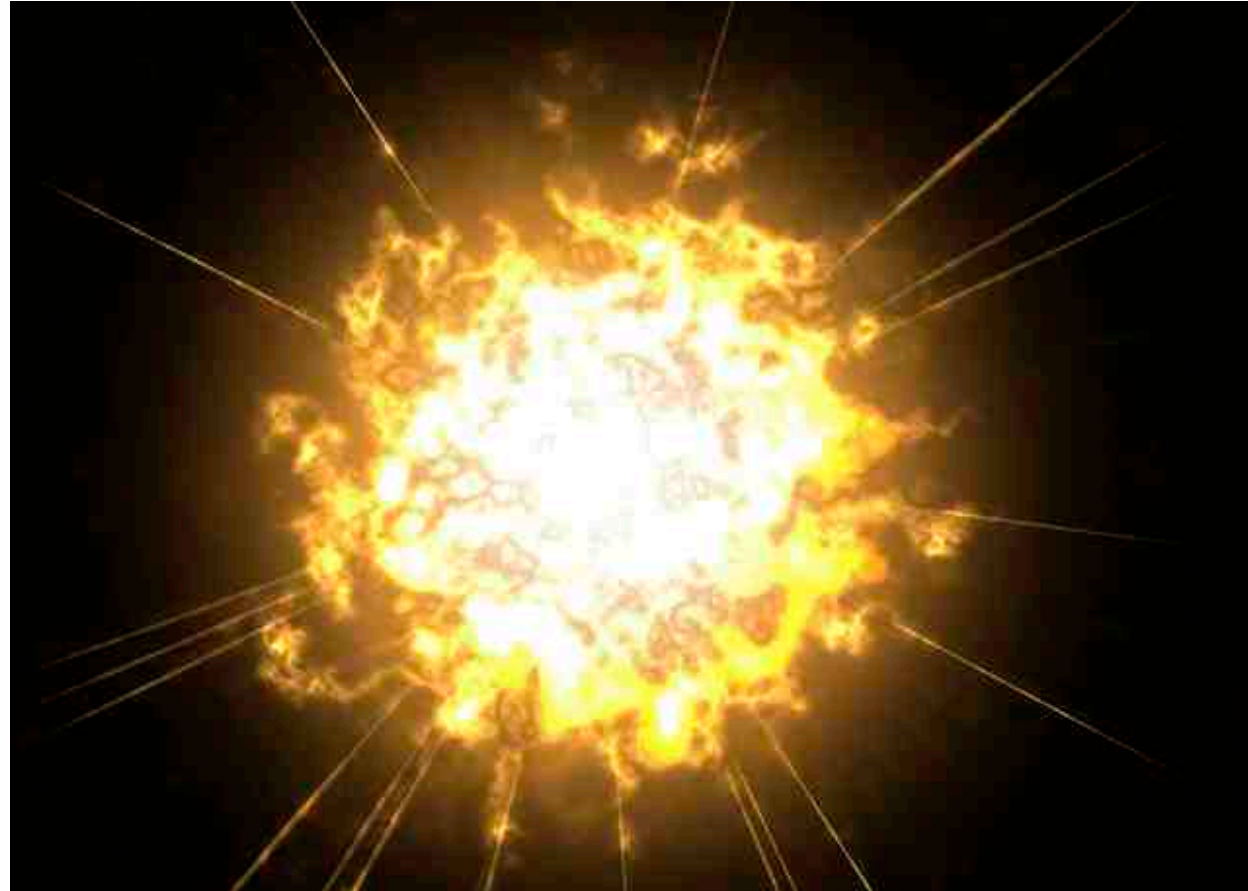




# Explosion des groupes



# Complication #1 : Big Bang



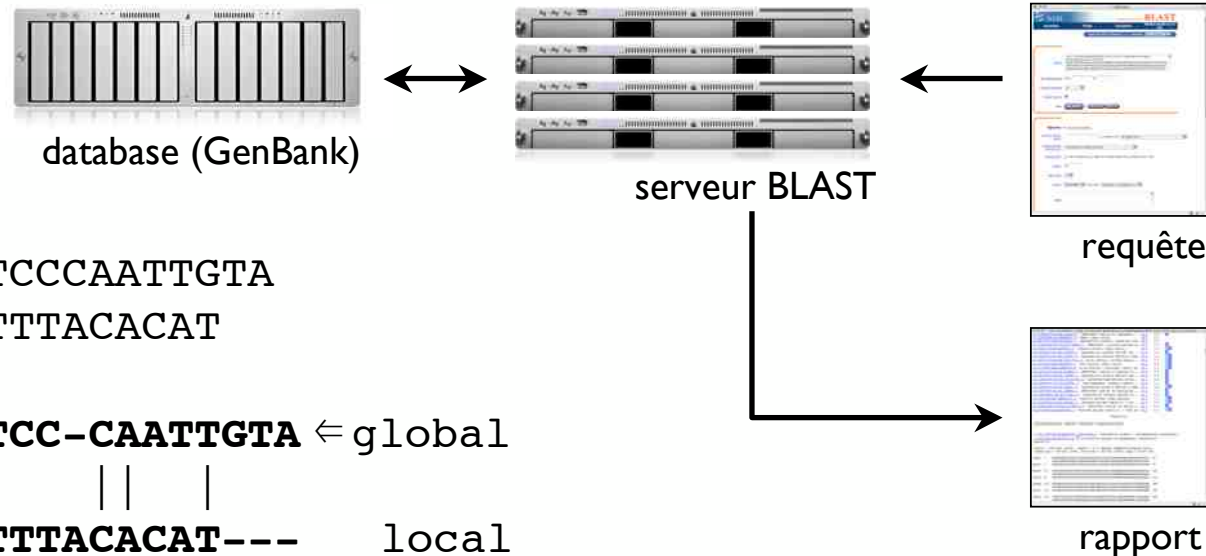
possible radiation **explosive** des Eucaryotes  
aggravée par le **signal non-phylogénétique**



*Transfert de gènes*  
*Orthologie, paralogie et xénologie*

# Basic Local Alignment Search Tool

- calcule la similarité entre deux séquences biologiques
- produit des alignements locaux : seule une portion de chaque séquence est alignée
- utilise des statistiques sophistiquées pour déterminer si un alignement pourrait avoir été obtenu par le seul hasard



```

S1 TTGACACCCTCCCAATTGTA
S2 ACCCCAGGCTTTACACAT

S1 TTGACACCCTCC-CAATTGTA ← global
   ||  ||  ||  |
S2 ACCCAGGCTTTACACAT--- local
                        ↓
S1 -----TTGACACCCTCCCAATTGTA
           ||  ||  ||
S2 ACCCCAGGCTTTACACAT-----
    
```

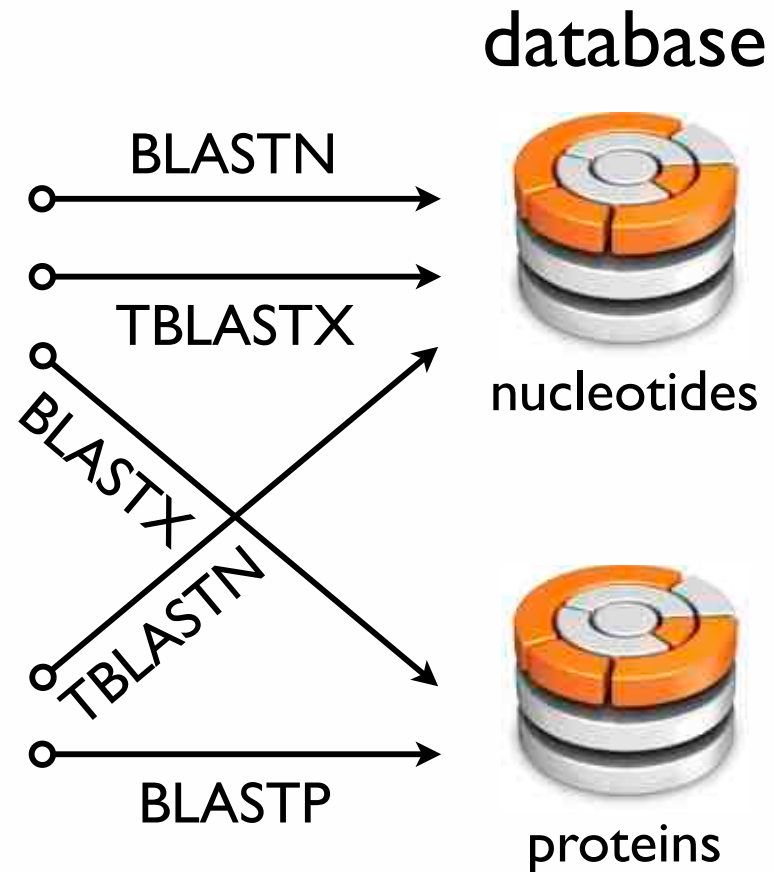
# BLAST

à chaque recherche sa variante de BLAST

query

```
>DNA_seq  
TATGGCAATTTAAAATTGGTATCAATGGTTTTGG  
TCGTATCGGCCGTATCGTATTCGGTGCAGCACA  
ACACCGTGATGACATTGAAGTTGTAGGTATTAA  
CGACTTAATCGACGTTGA
```

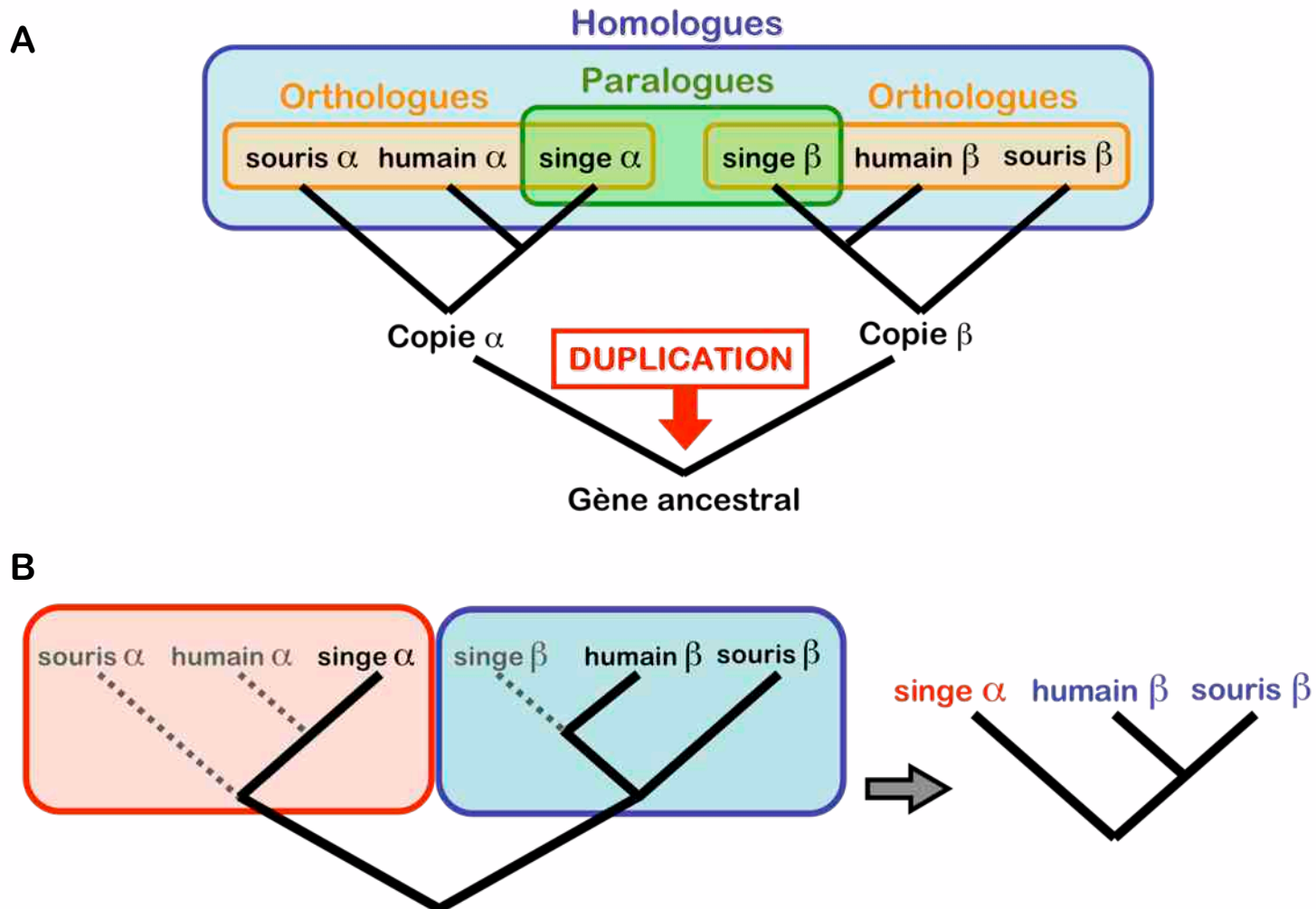
```
>protein_seq  
MLQTAPMLPGLGPHLVPQLGALASASRLLGSIA  
SVPPQHGGAGFQAVRGFATGAVSTPAASSPGHK  
PAATHAPPTRLDLKPGAGSFAAGAVAPHPGINP  
ARMAADSASAA
```



Les recherches de type protéique sont **plus sensibles**.

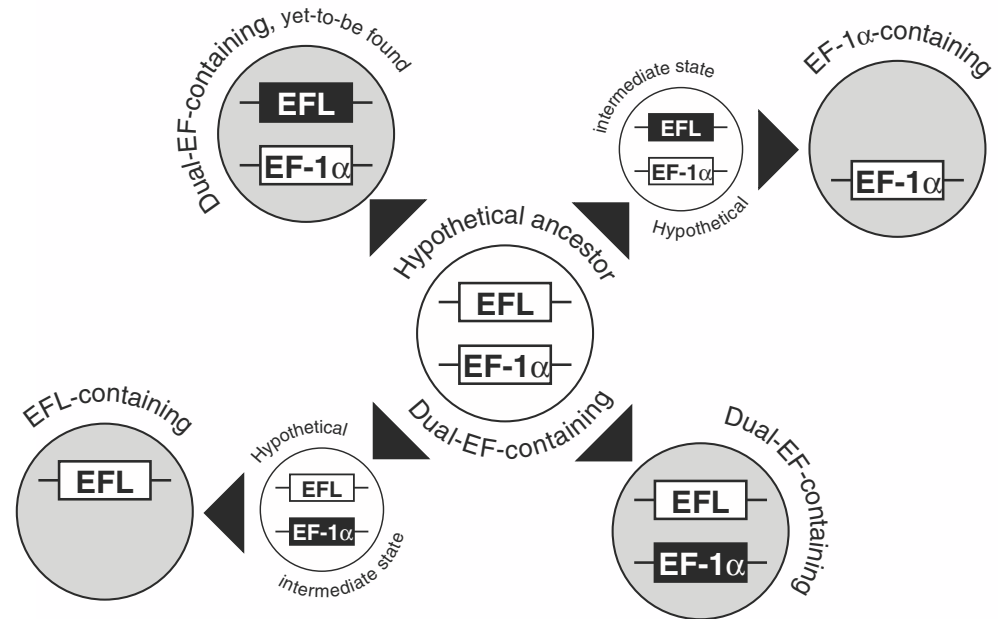
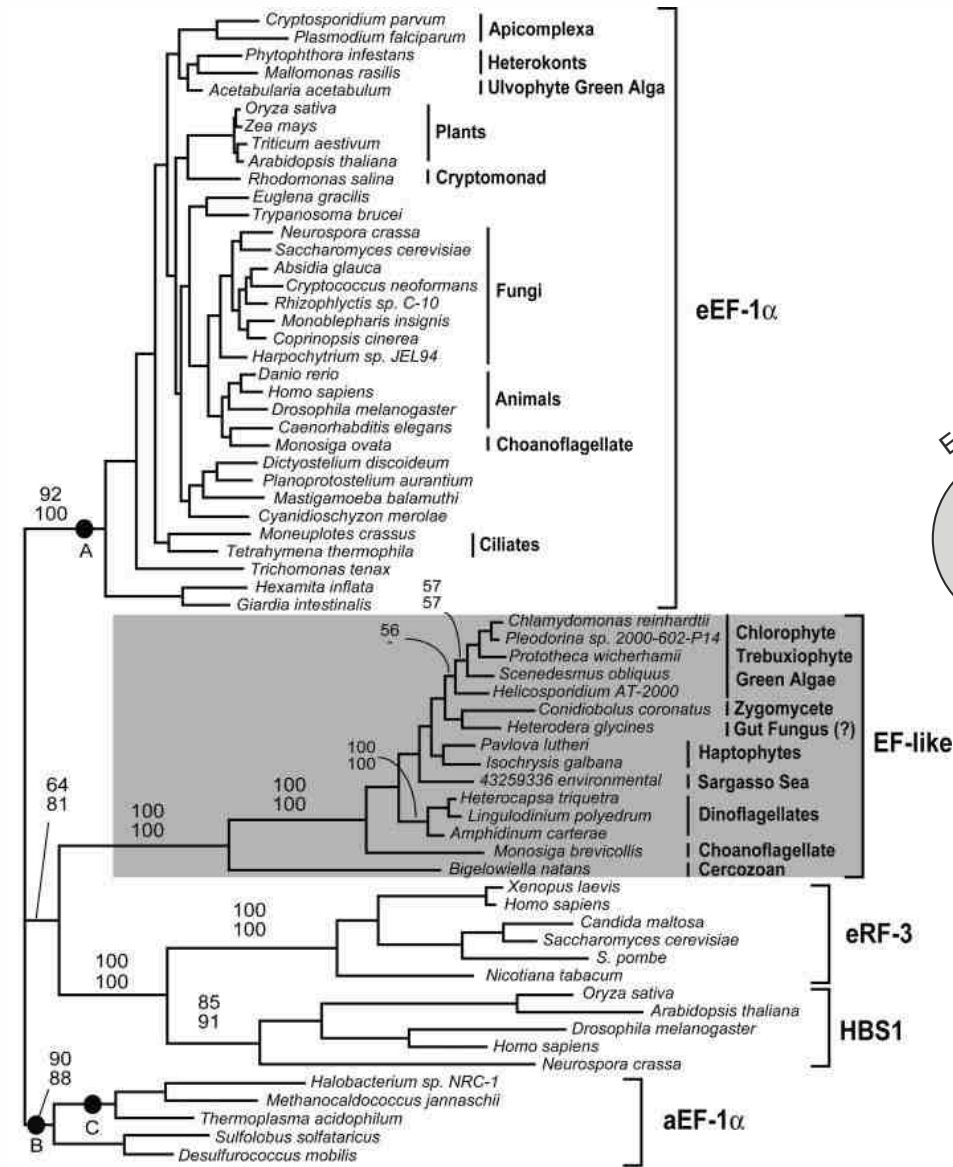
# Orthologie et Paralogie

## les dangers de la paralogie cachée





# Orthologie et Paralogie

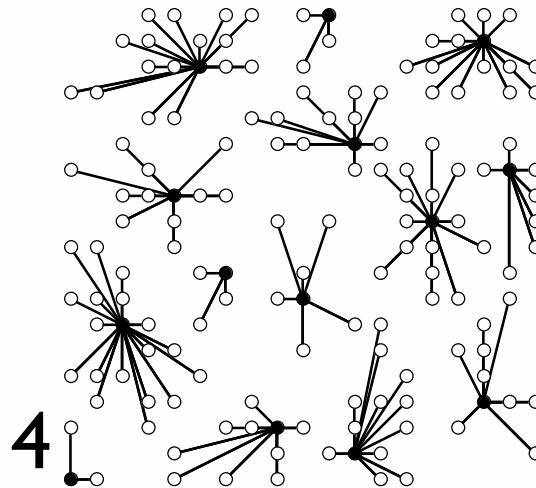
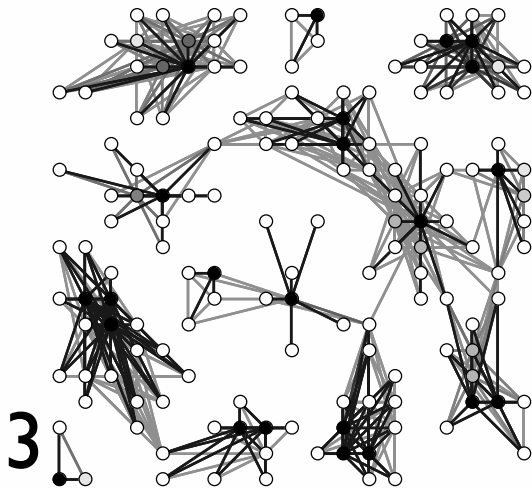
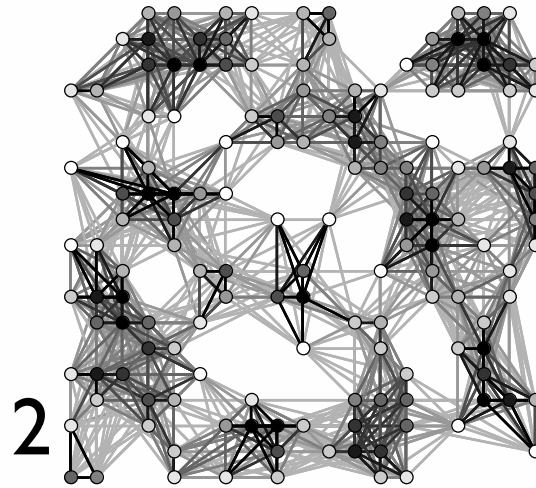
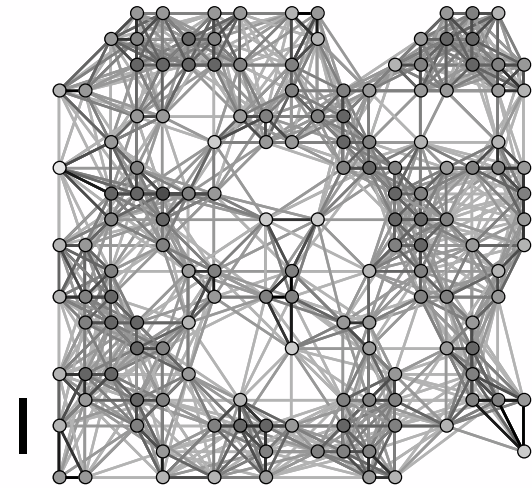


un exemple de **paralogie cachée** profonde :  
le couple **eEF-1a / EF-like**



# Orthologie et Paralogie

détermination des groupes orthologues



OrthoMCL  
objective le  
processus, mais  
n'est pas sans faille.

# Endosymbioses primaires

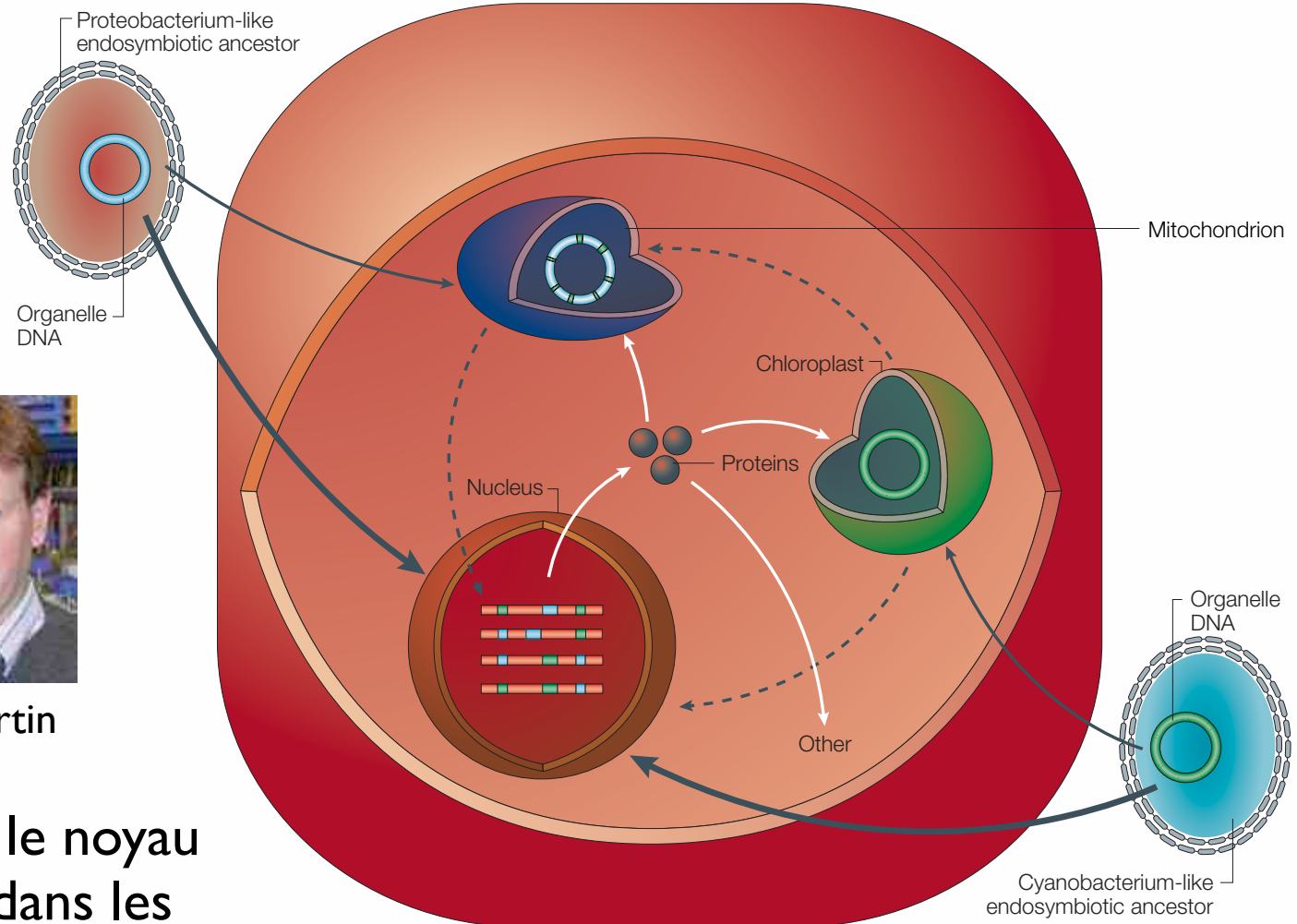


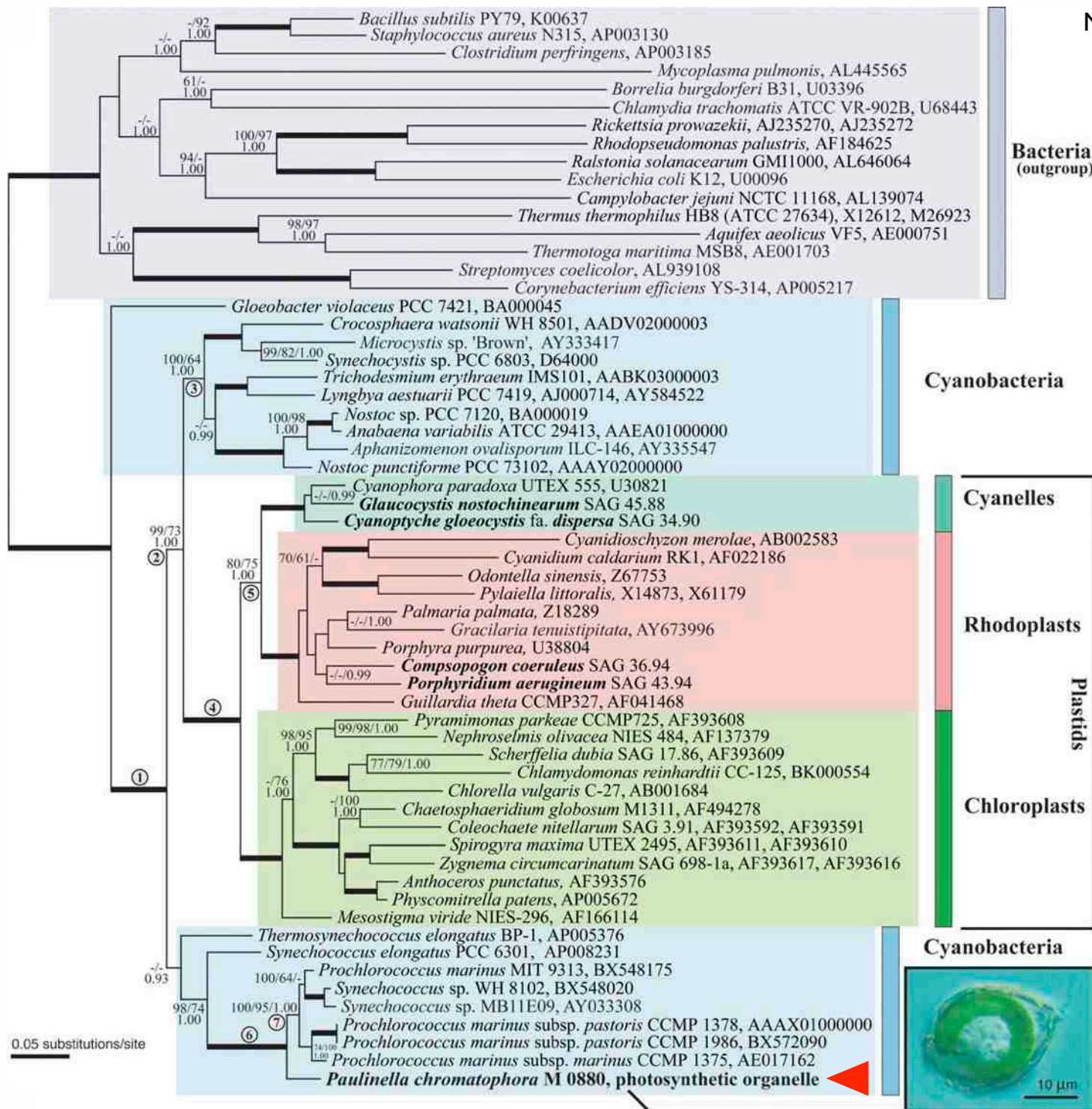
Lynn Margulis



Bill Martin

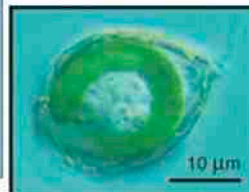
relocalisation dans le noyau  
de gènes encodés dans les  
organites endosymbiotiques



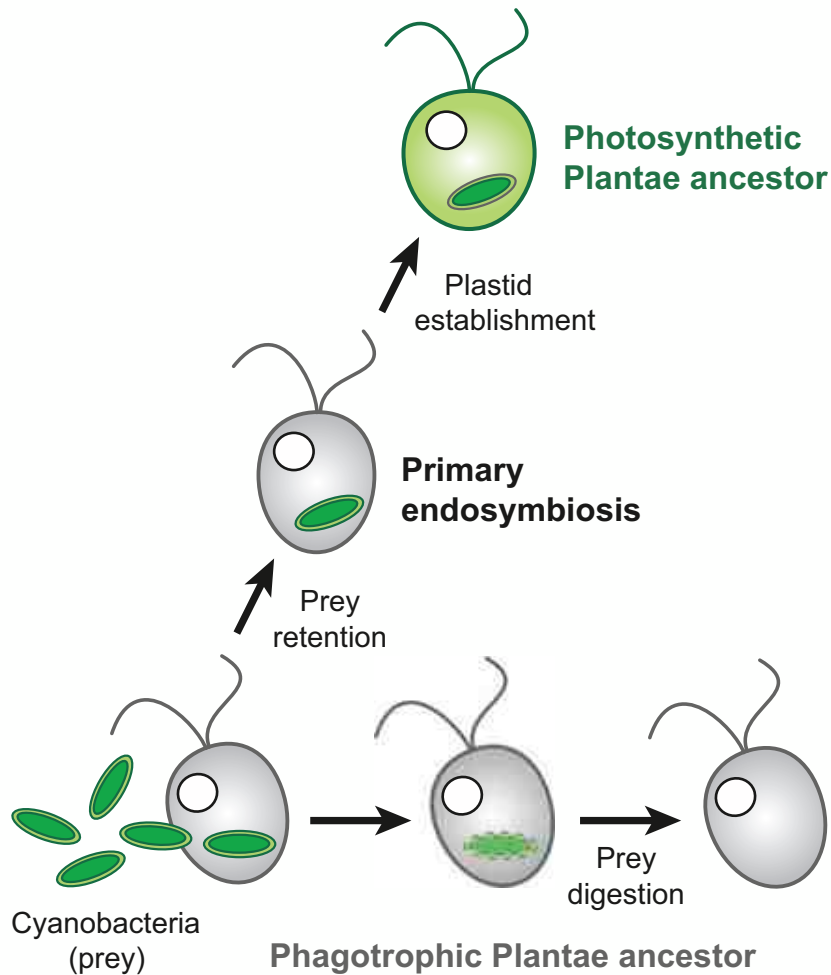


opéron rRNA

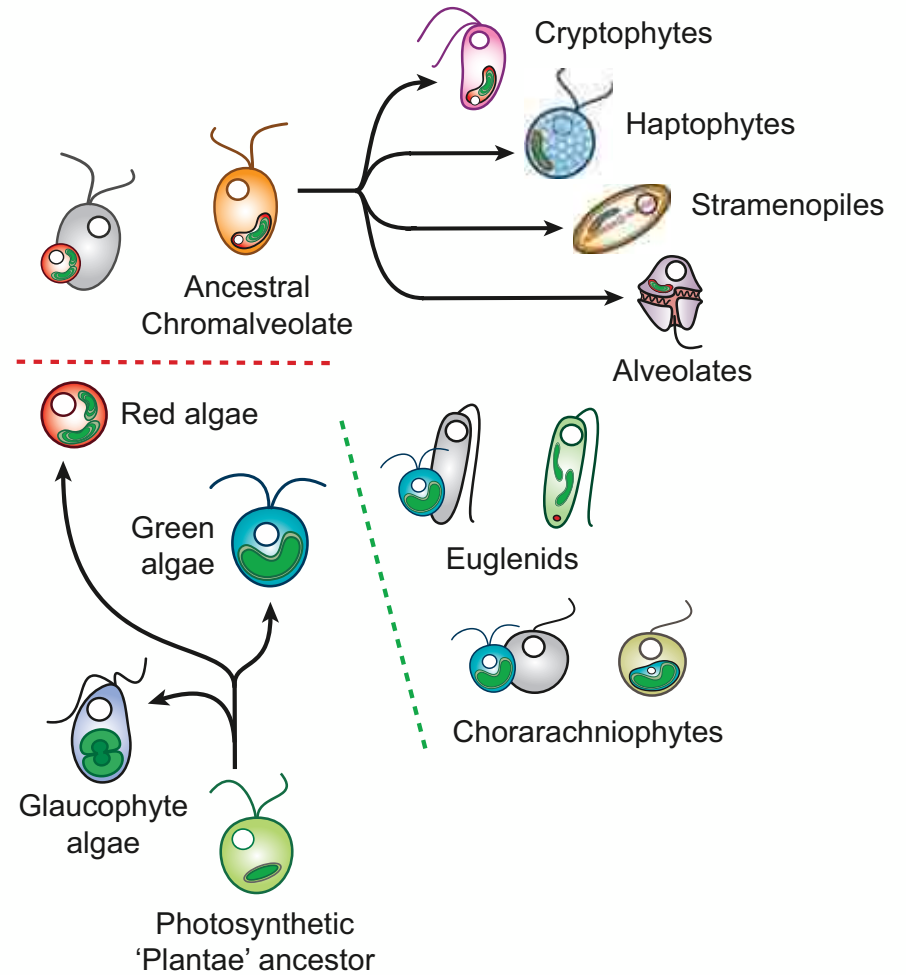
une seconde  
endosymbiose  
primaire :  
*Paulinella*



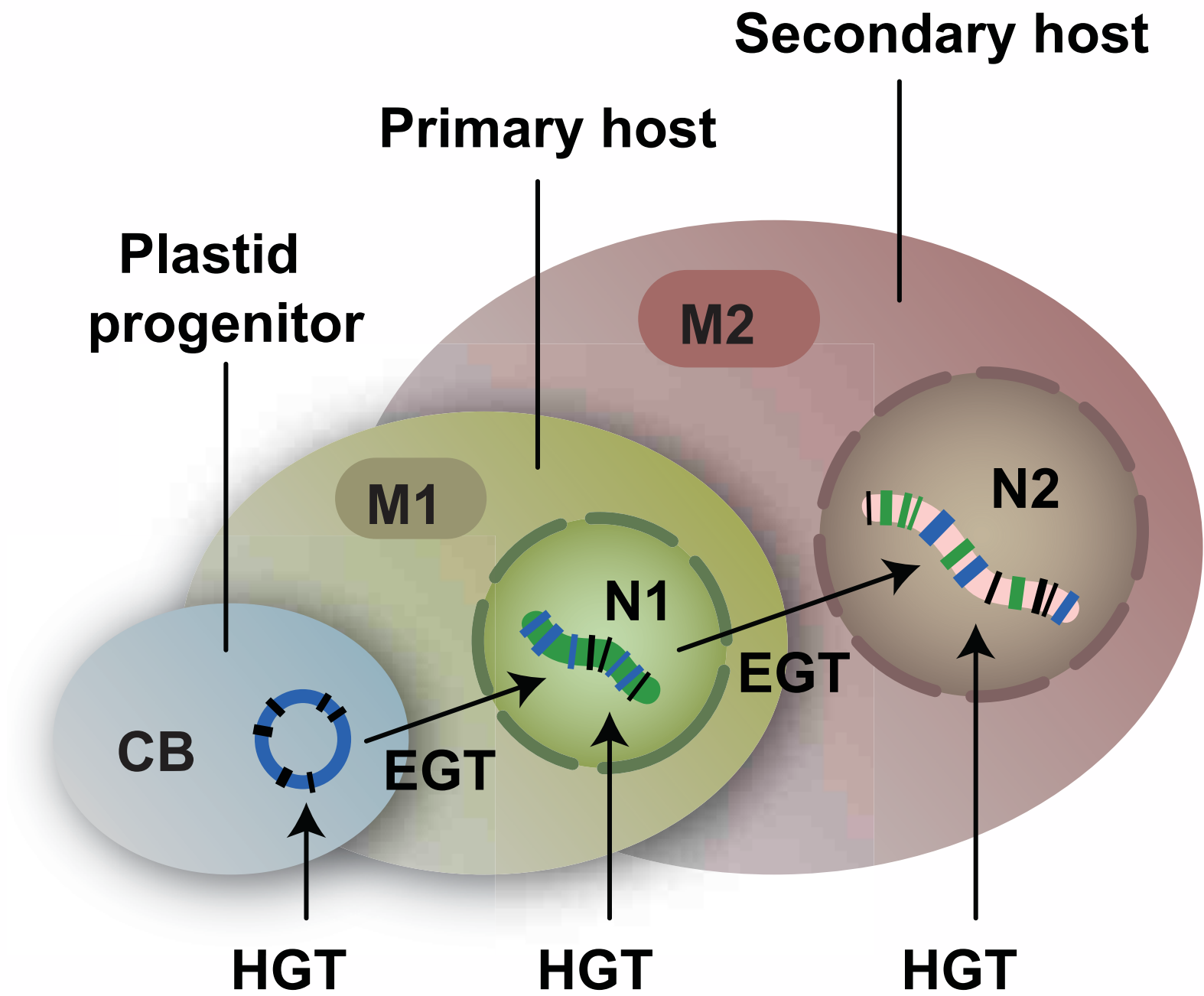
## endosymbiose primaire Plantae

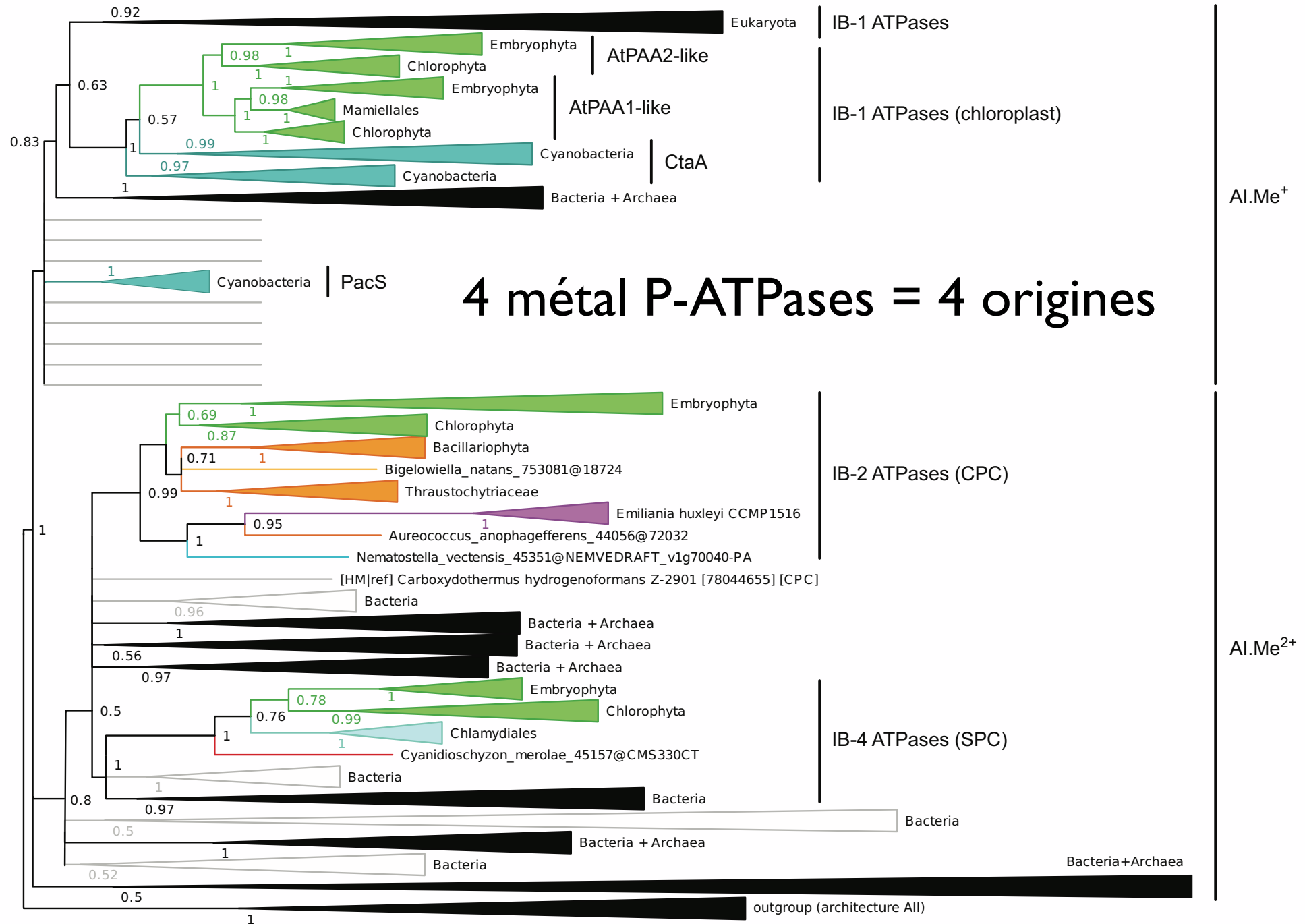


## endosymbiose secondaire « Chromalveolés »



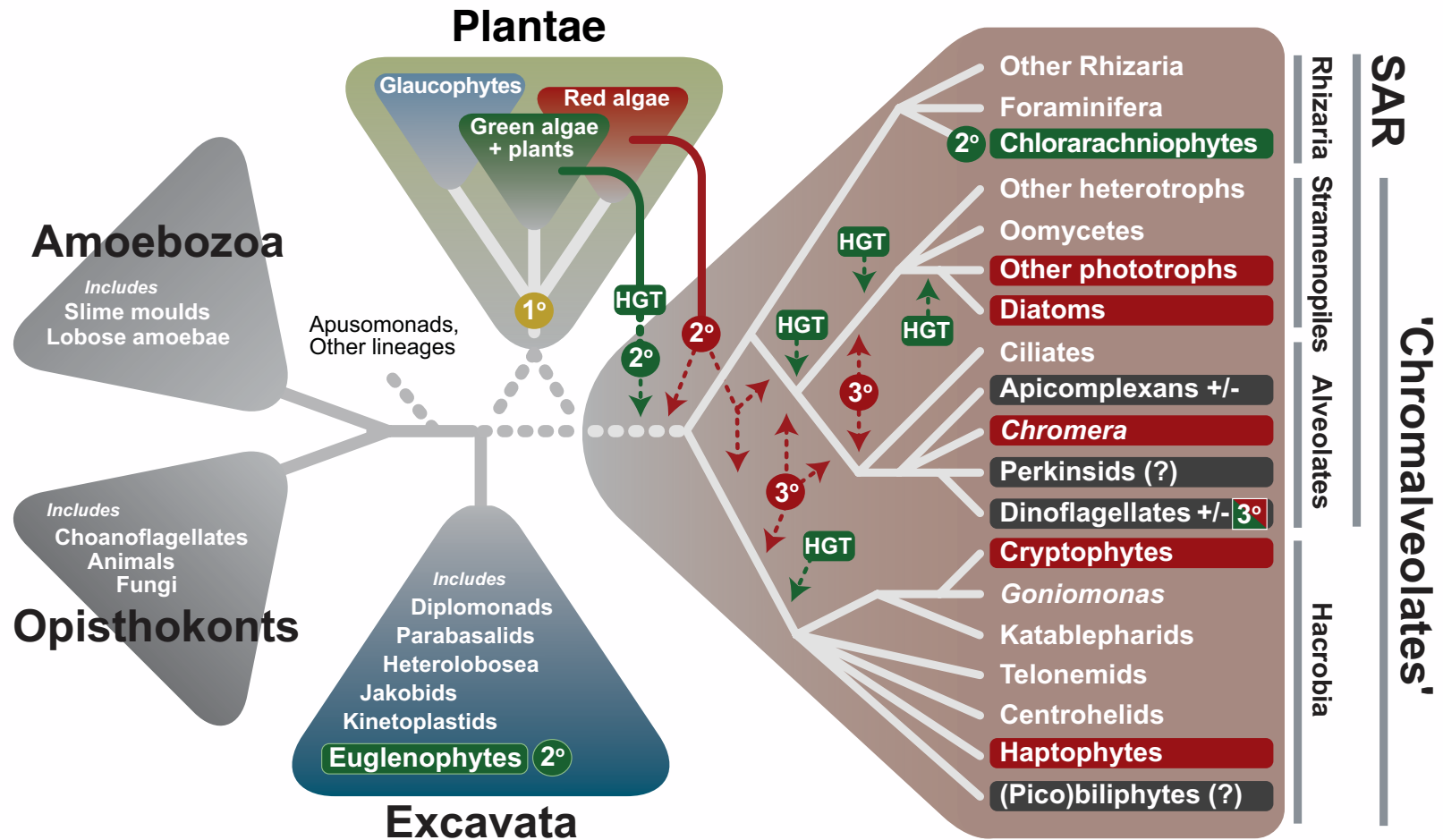






0.4

# Complication #2 : EGT

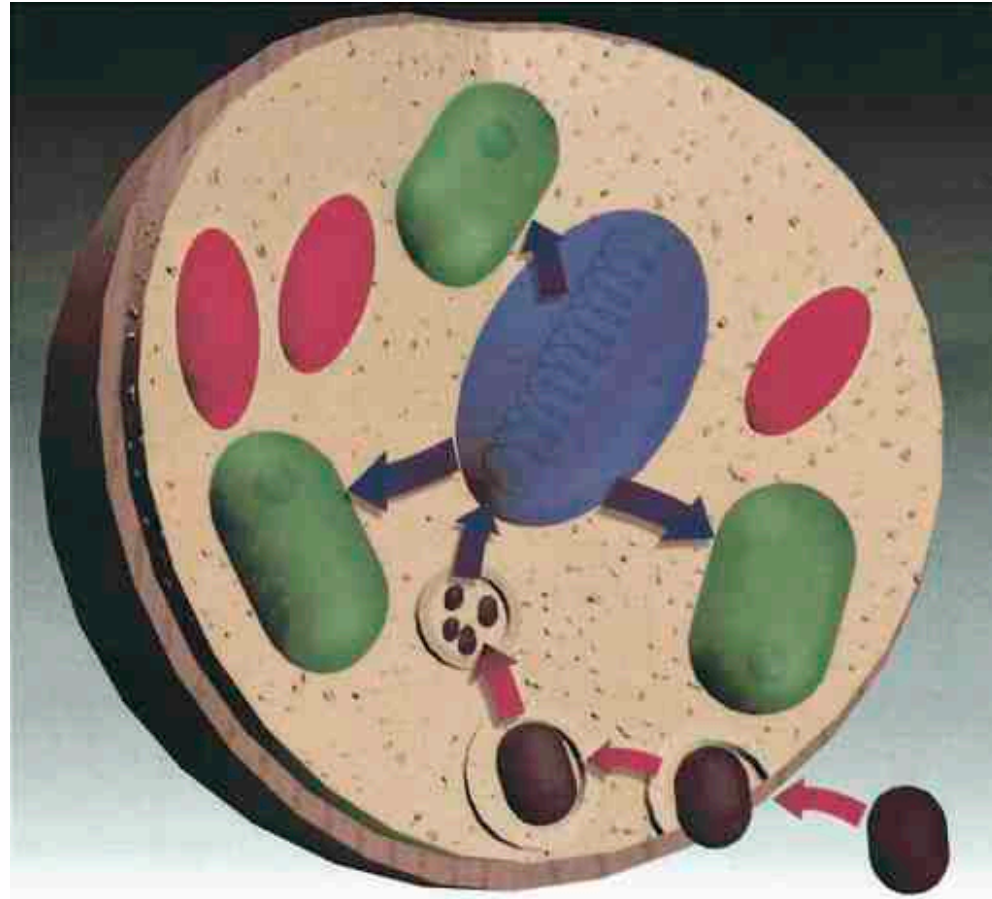


chimérisme complexe (mais biaisé) des génomes algaux  
lié aux multiples transferts indépendants de plastes

# Complication #3 : HGT



W. Ford Doolittle

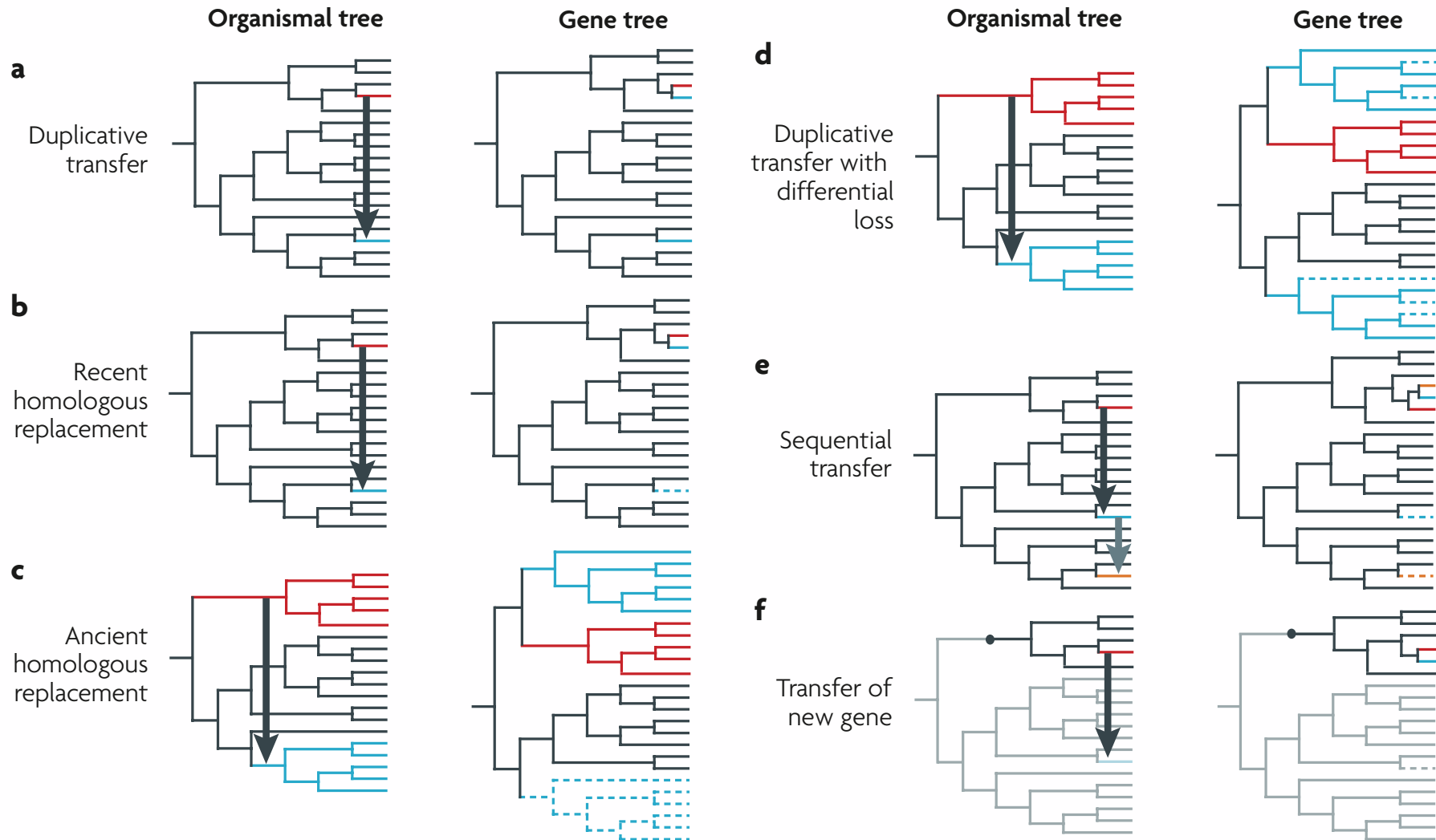


transfert **horizontal** de gènes  
lié au « You are what you eat »



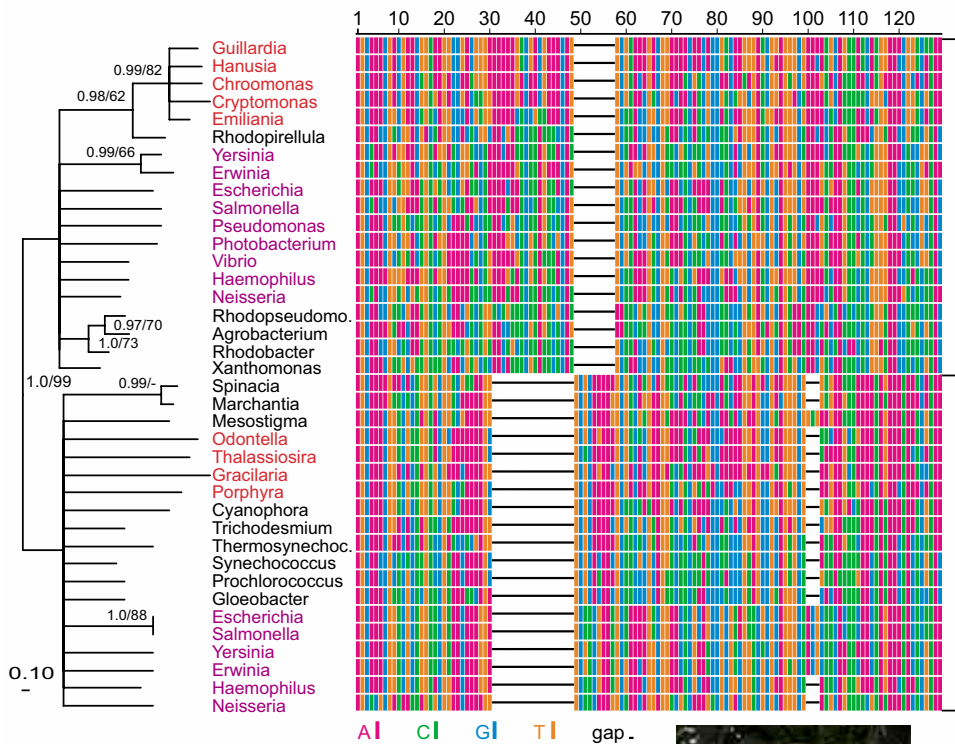
# Transfert de gènes

évolution des gènes  $\neq$  évolution des organismes



# Transfert de gènes

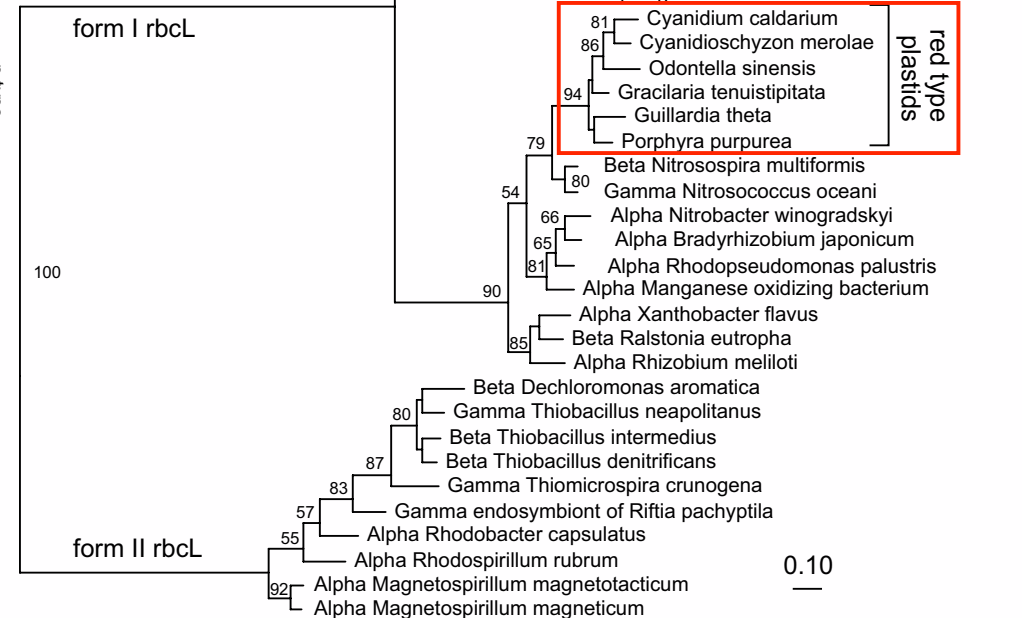
## aussi dans les génomes des organites



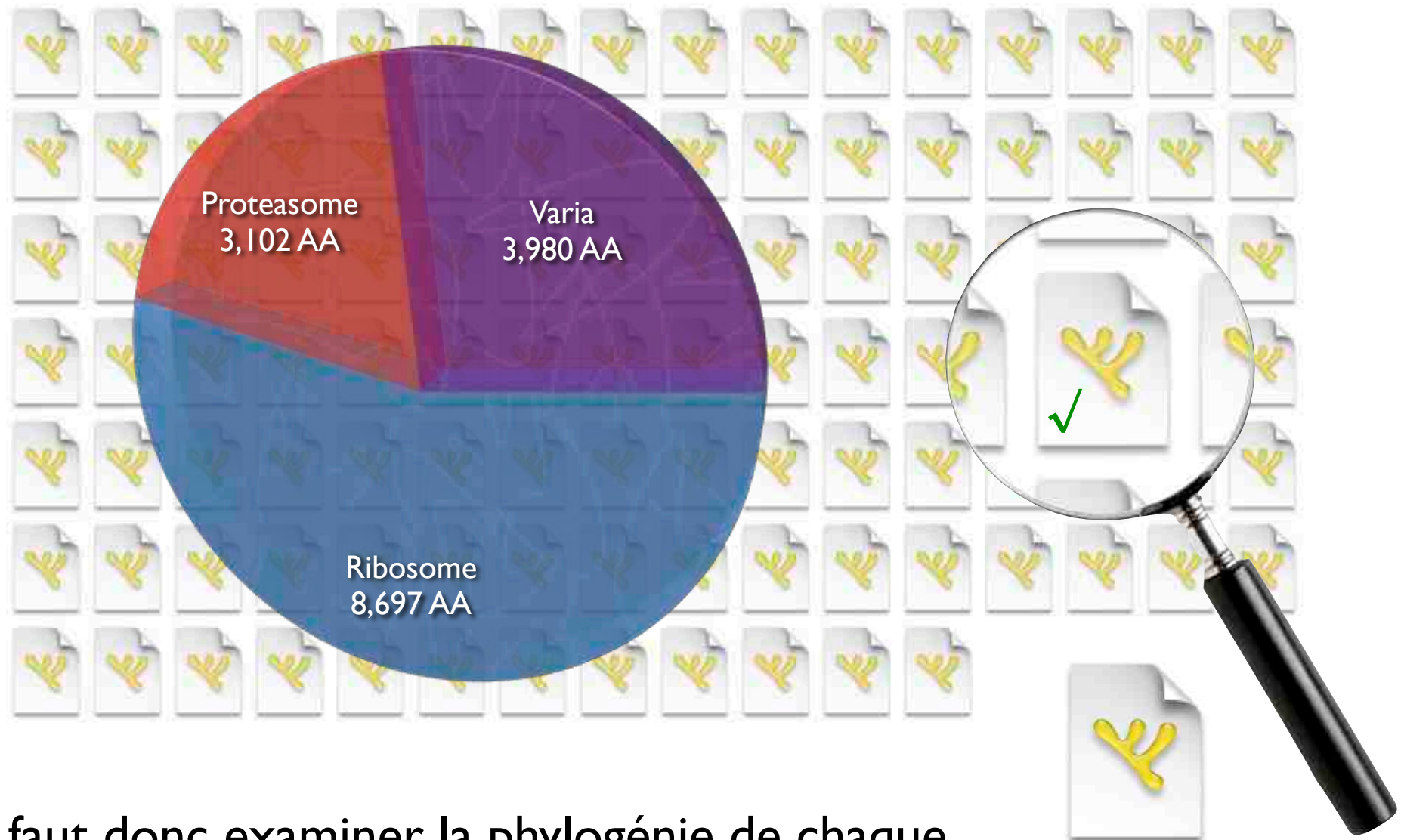
Jeff Palmer

**rps36**

**rbcl**

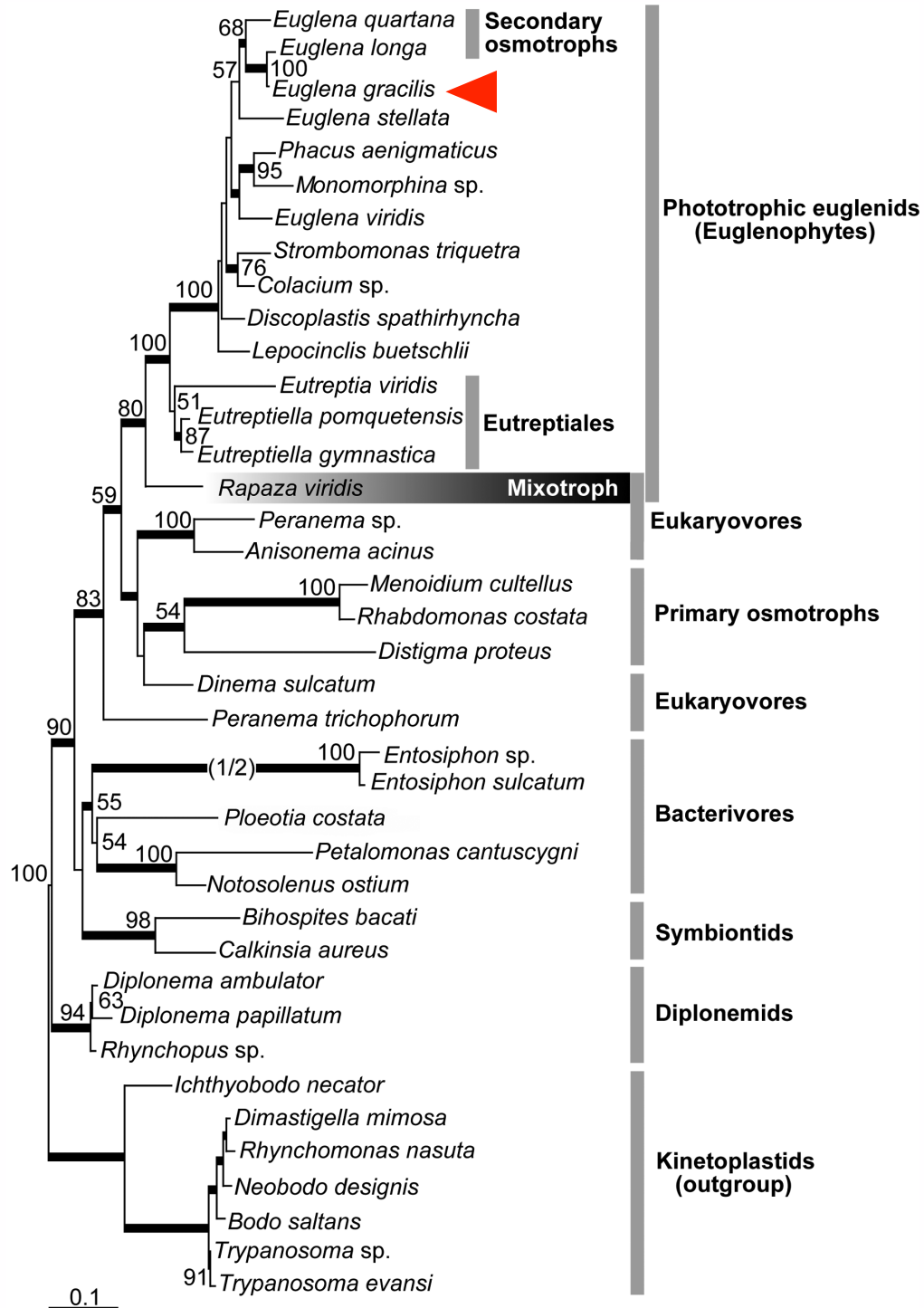


0.10



arbre  
concaténé

Il faut donc examiner la phylogénie de chaque gène et la comparer à celle obtenue à partir de la supermatrice (analyse de **congruence**).

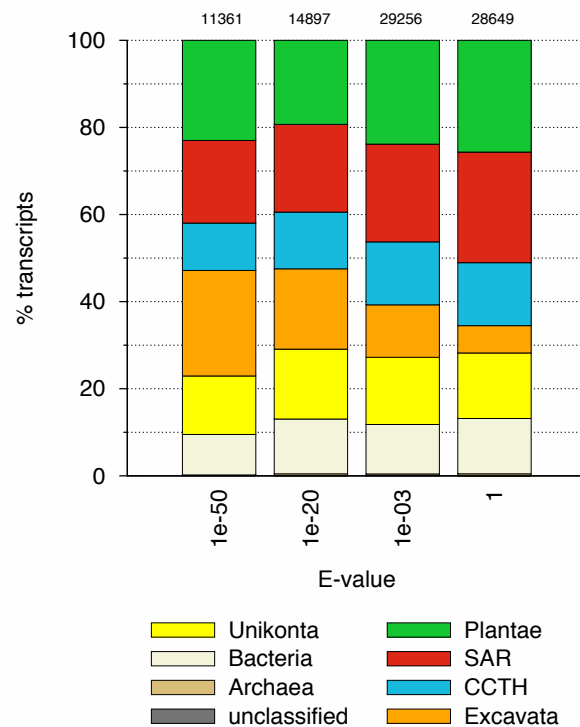


une algue « verte » au passé trouble : *Euglena gracilis* (Discoba, Excavata)

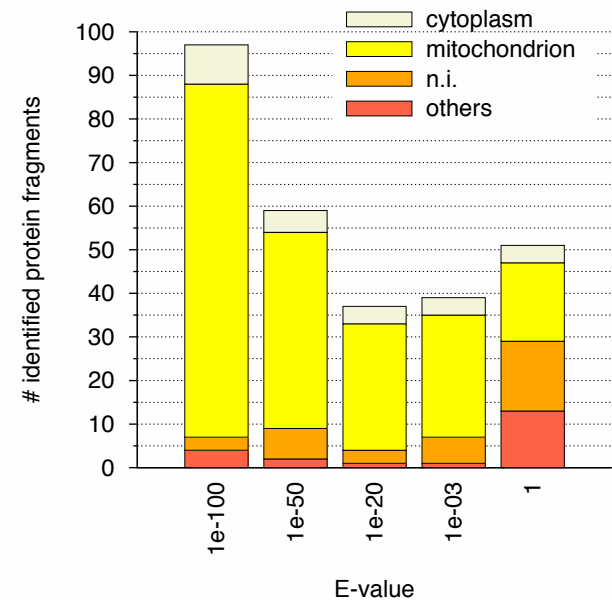
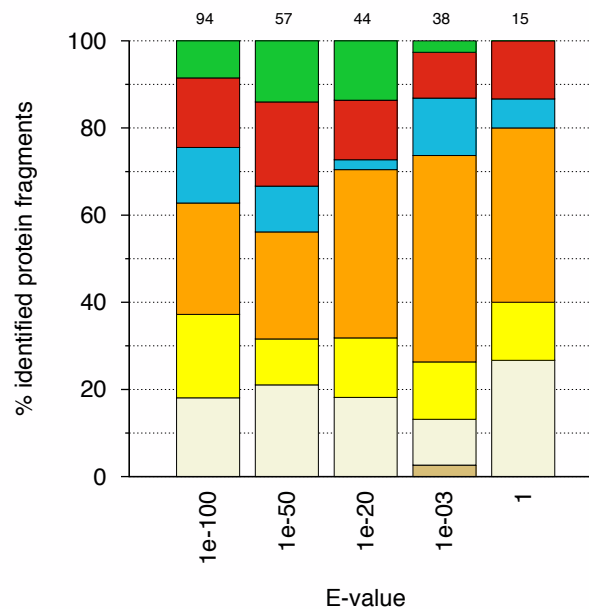
# Transfert de gènes

origines diverses des gènes chez *Euglena gracilis*

taxonomie  
des transcrits



taxonomie et localisation  
des protéines mitochondriales



combinaison EGT/HGT



Diaphoretickes

ARCHAEPLASTIDA

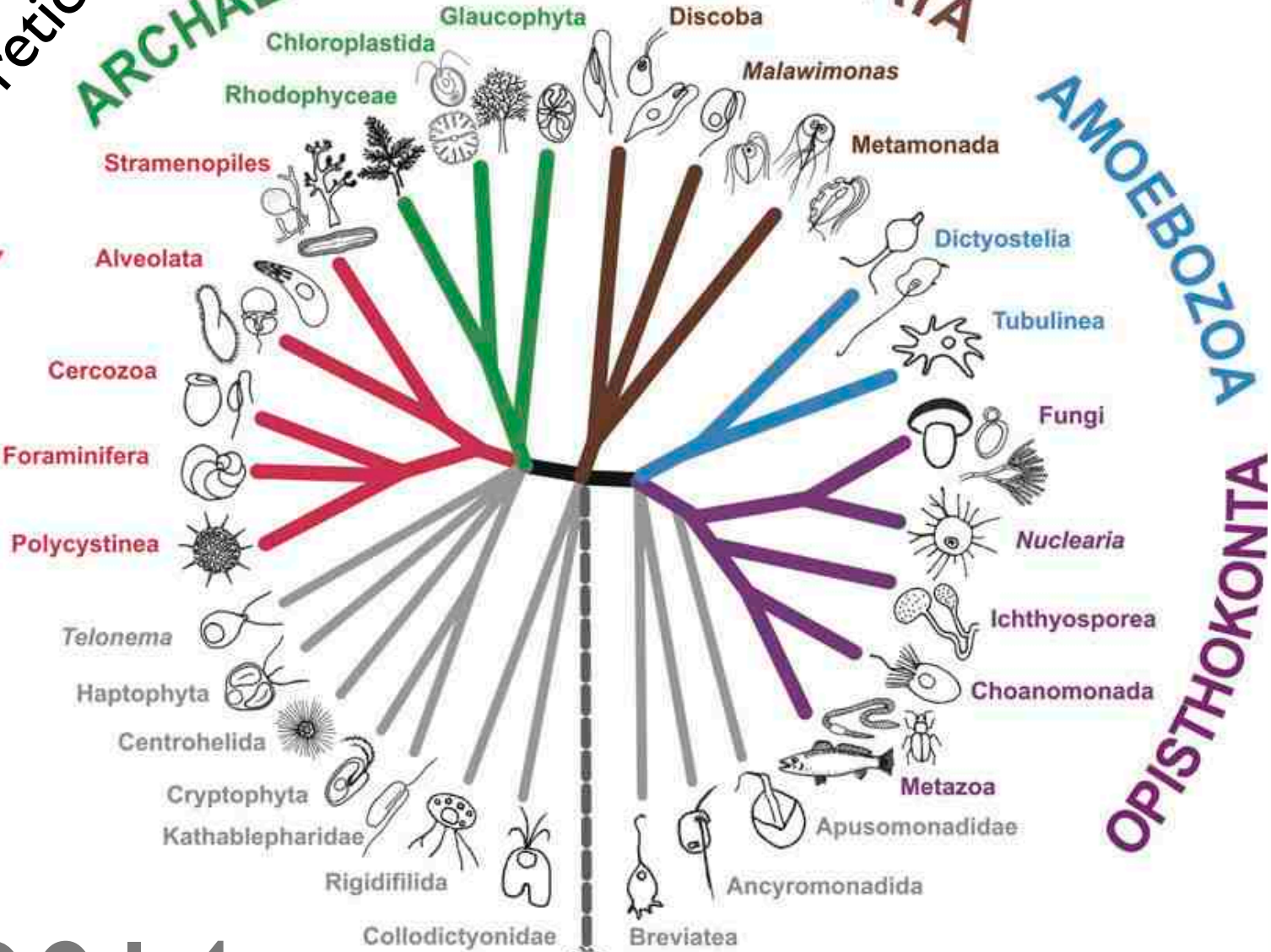
EXCAVATA

AMOEBOZOEA

OPISTHOKONTA

Amorphea

SAR

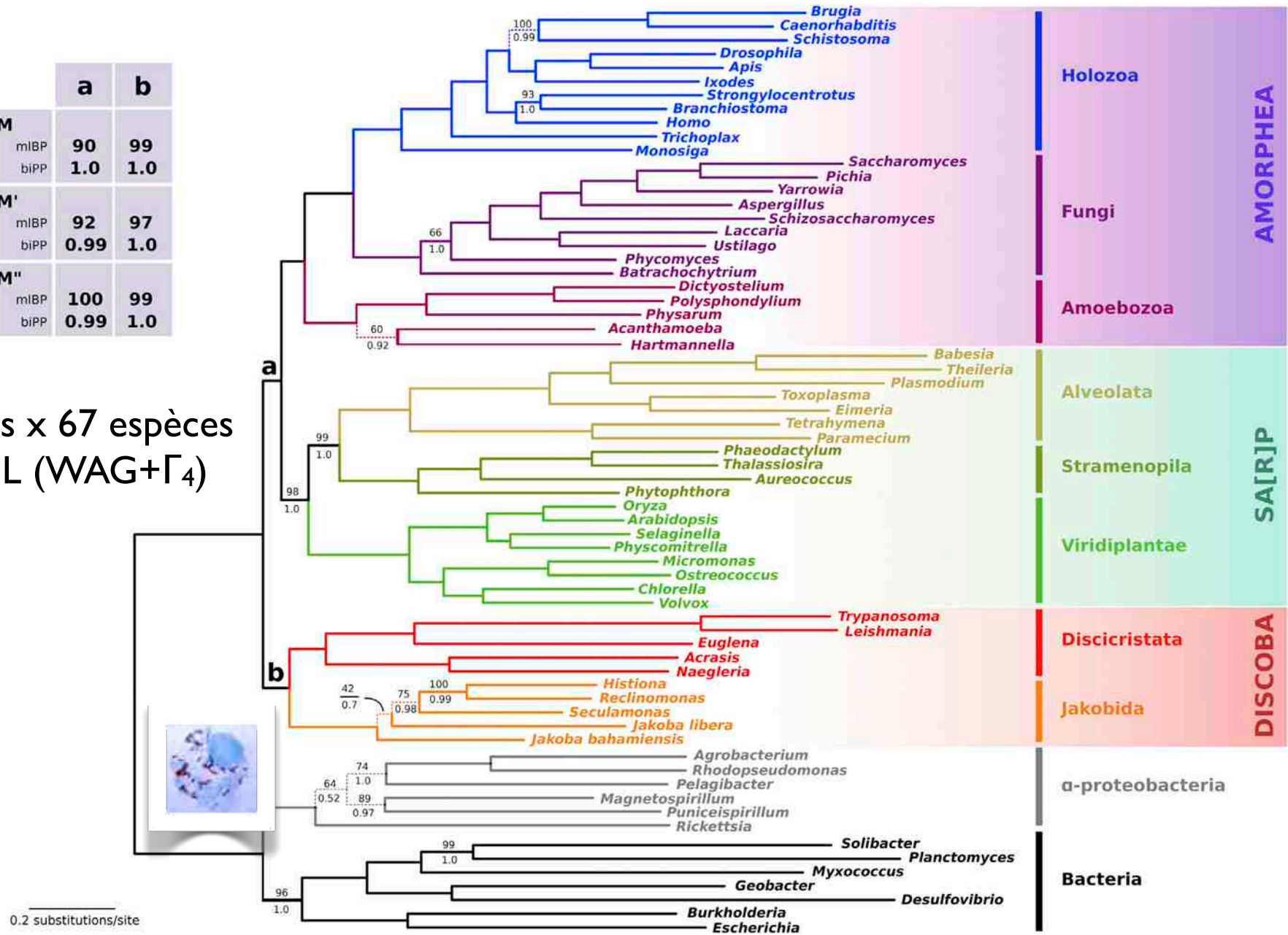


~2014

ARCHAEBACTERIA EUBACTERIA

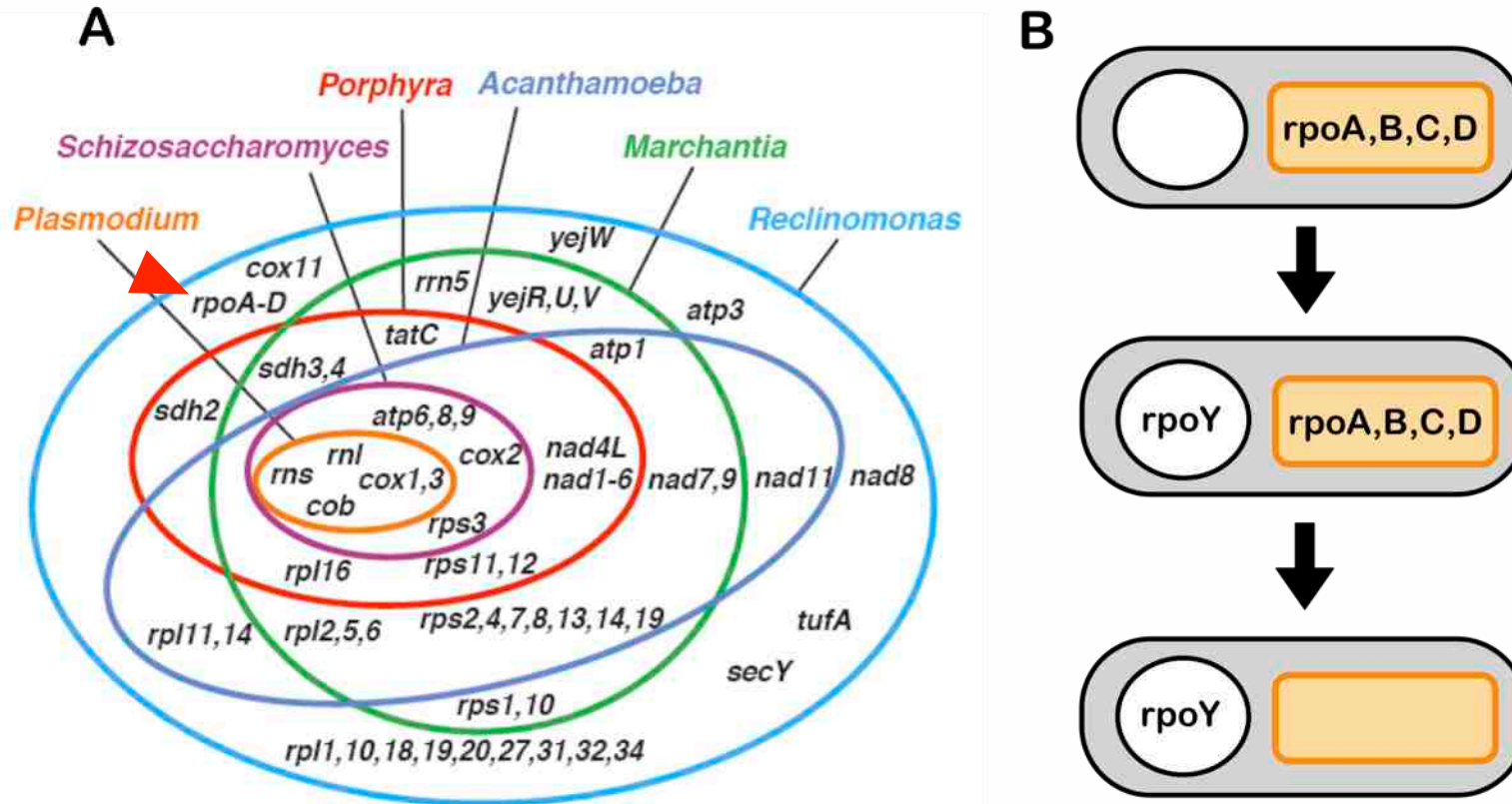
	a	b
<b>M</b>		
mIBP	90	99
biPP	1.0	1.0
<b>M'</b>		
mIBP	92	97
biPP	0.99	1.0
<b>M''</b>		
mIBP	100	99
biPP	0.99	1.0

37 gènes x 67 espèces  
RAxML (WAG+ $\Gamma_4$ )



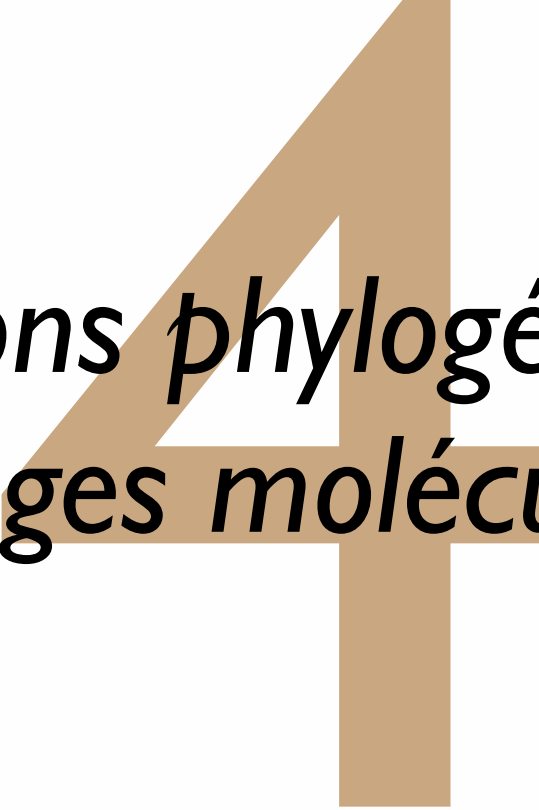
position de la racine des eucaryotes dans les Excavata ?

# Discoba : Jakobida



Les jakobidés sont les seuls eucaryotes avec une **ARN polymérase alpha-protéobactérienne**. Tous les autres ont une enzyme (nucléaire) transférée depuis un bactériophage.

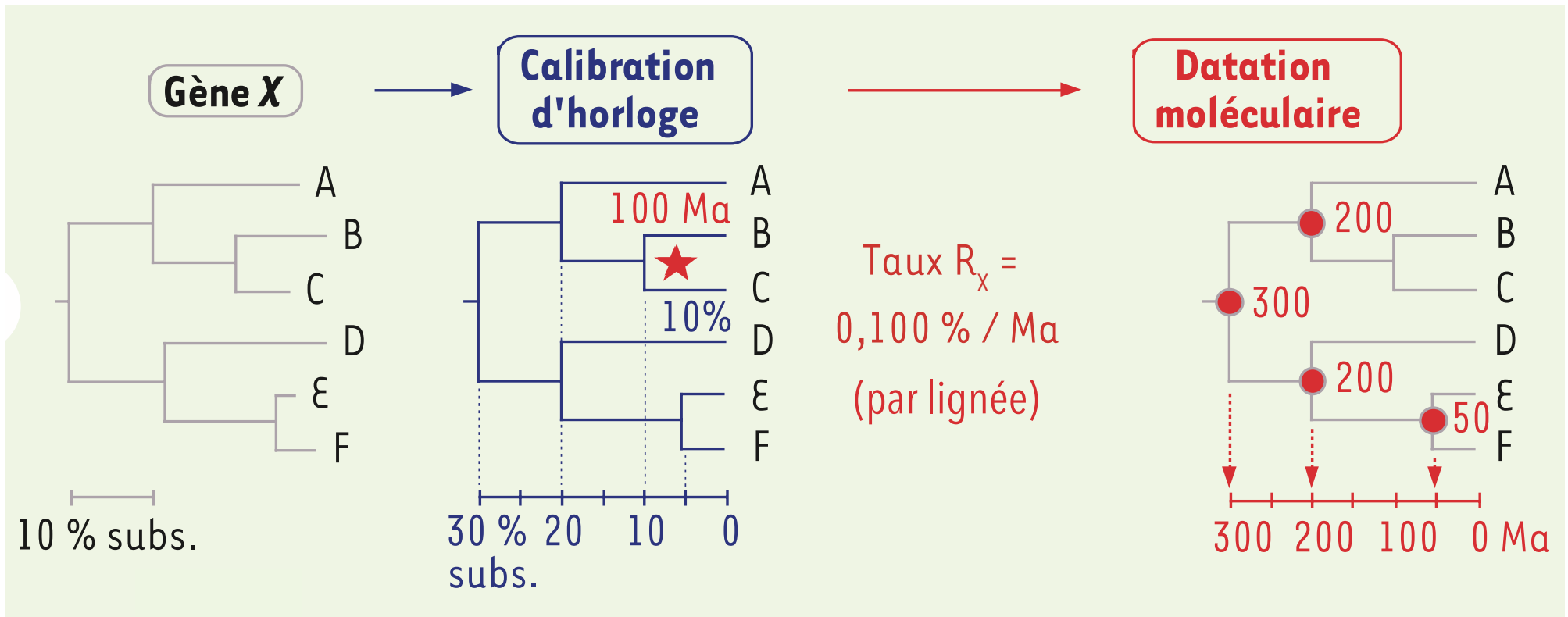




*Prédictions phylogénétiques*  
*Horloges moléculaires*

# Horloges moléculaires

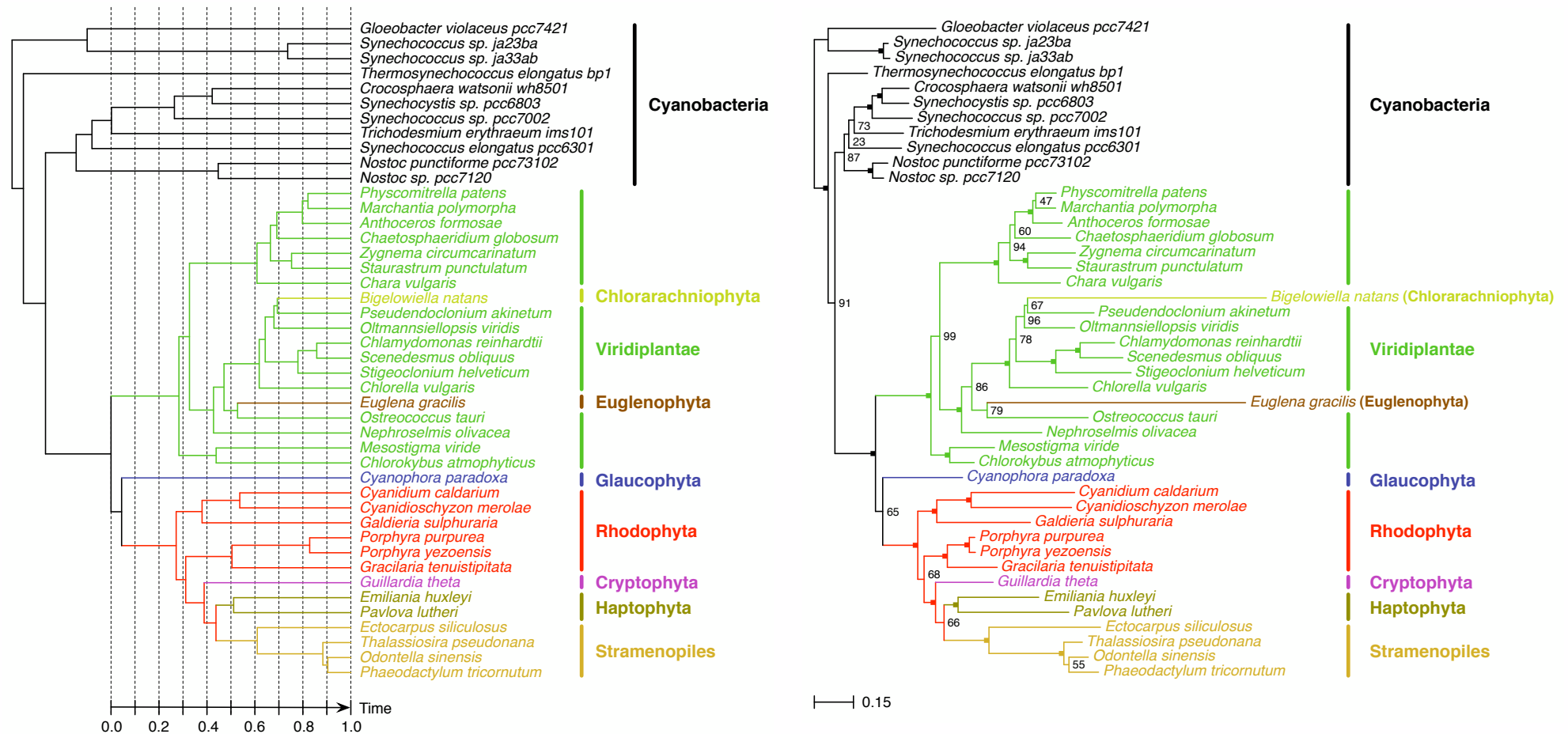
principe de base (vitesse d'évolution constante)



Si les différences de vitesses entre branches ne sont pas importantes, on peut contraindre l'arbre à être **ultramétrique**. Cela permet d'inférer la **vitesse d'évolution** à partir d'un **seul point de calibration** et de dater tous les autres noeuds en conséquence.

# Horloges moléculaires

Hélas, les lignées n'évoluent pas à vitesse constante.



55 gènes x 44 espèces  
PhyloBayes (CAT+ $\Gamma_4$ )



distance = temps x vitesse

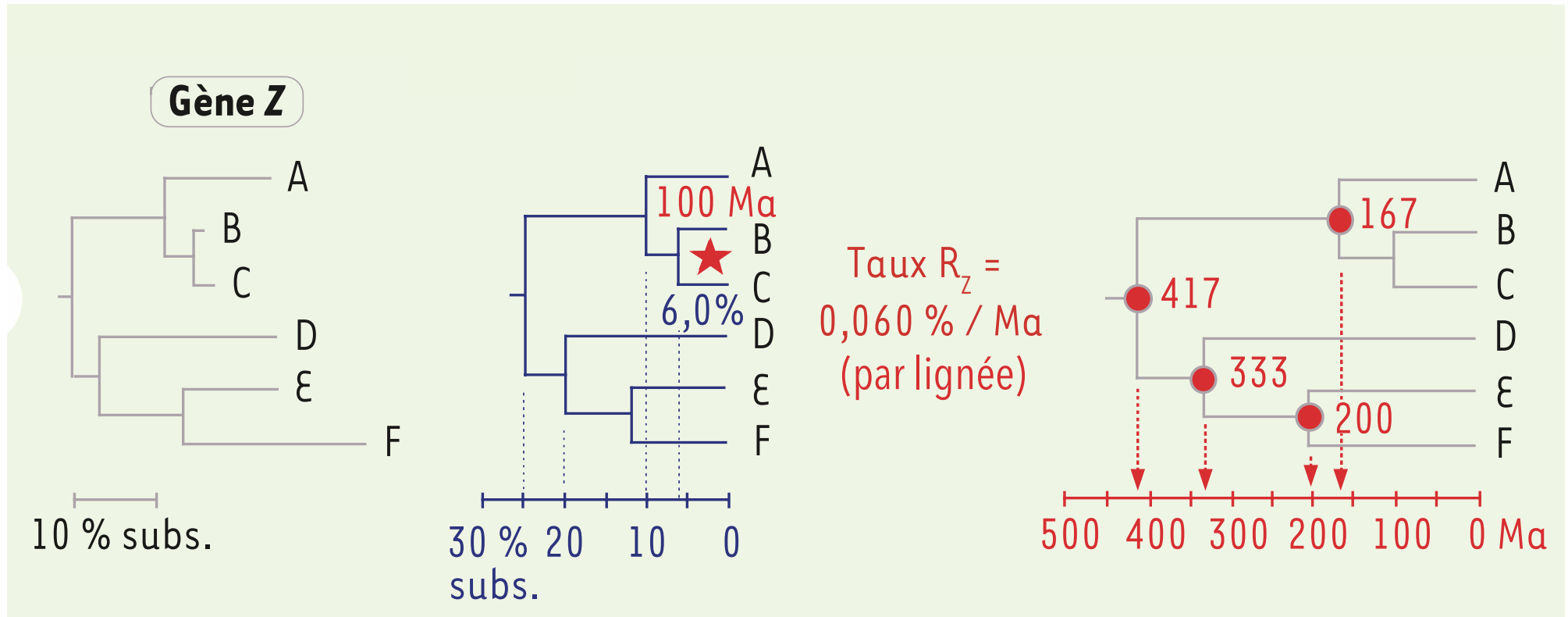




distance = temps x vitesse

# Horloges moléculaires

impact de la vitesse d'évolution sur la datation

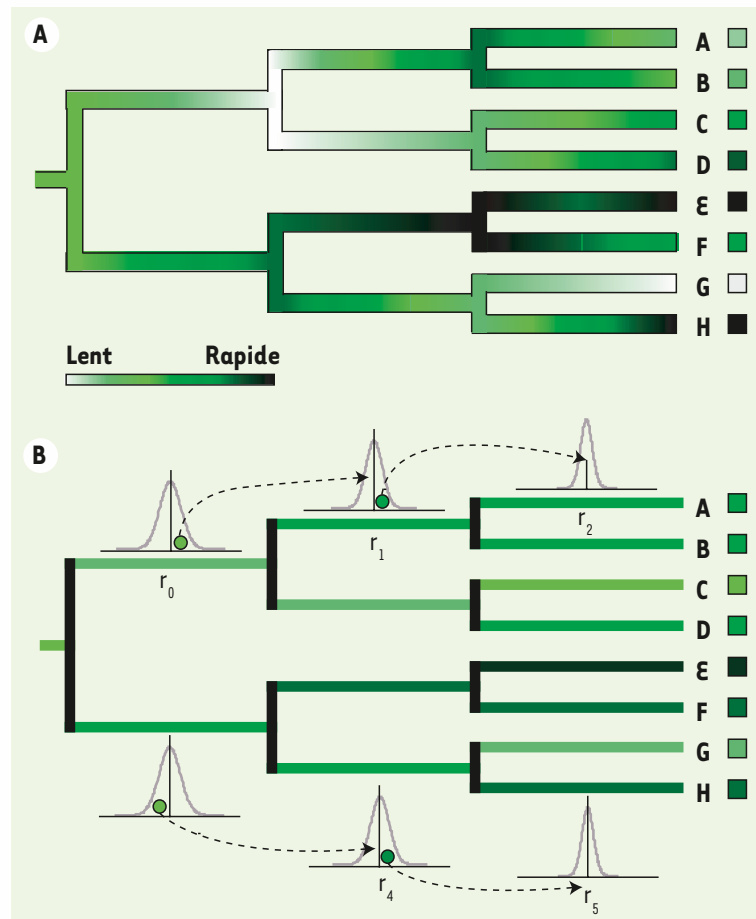


Si nous utilisons le gène Z, dont les vitesses d'évolution sont très différentes entre lignées (pas d'ultramétrie), les temps de divergence inférés seront incorrects puisque l'horloge n'est pas constante.

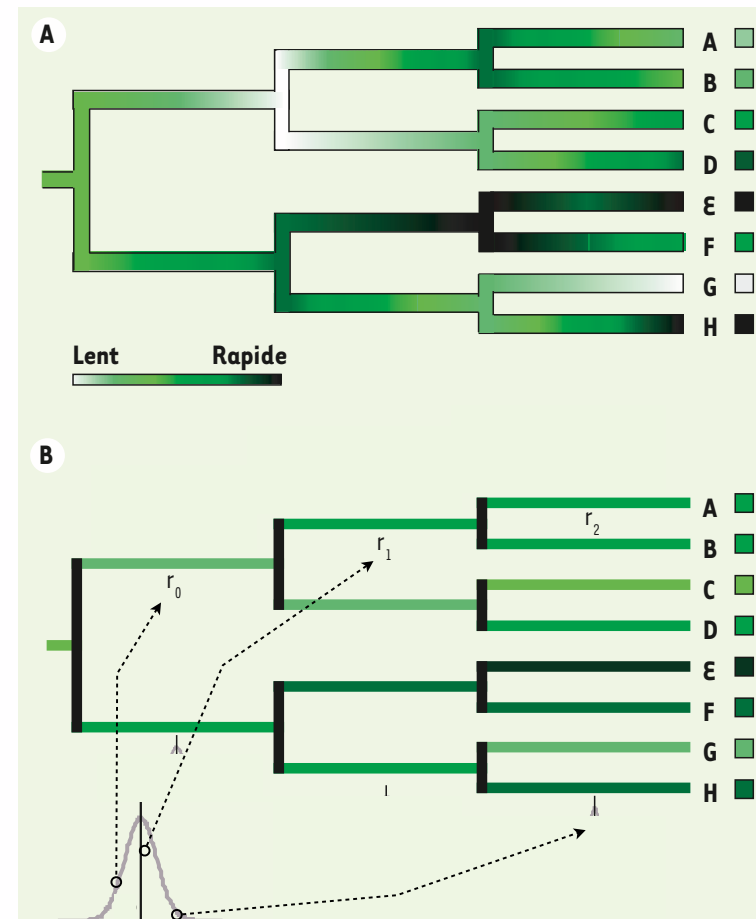
# Horloges moléculaires

solution : les horloges moléculaires assouplies

vitesse corrélées



vitesse non-corrélées







Taxon	Fossil	Eon*	Calibration <sup>†</sup>		Ref(s).
			Min	Dist	
Amniota	<i>Westlothania</i>	Phan	328.3	4, 3	(54)
Angiosperms	Oldest angio pollen	Phan	133.9	2, 10	(55)
Ascomycetes	<i>Paleopyrenomycites</i>	Phan	400	4, 50	(56)
Coccolithophores	Earliest Heterococcolith	Phan	203.6	2, 8	(57)
Diatoms	Earliest diatoms	Phan	133.9	2, 100	(58)
Dinoflagellates	Earliest gonyaulacales	Phan	240	2, 10	(59)
Embryophytes	Land plant spores	Phan	471	2, 20	(60)
Endopterygota	Mecoptera	Phan	284.4	5, 5	(61)
Eudicots	Eudicot pollen	Phan	125	2, 1.5	(62, 63)
Euglenids	<i>Moyeria</i>	Phan	450	2, 40	(64)
Foraminifera	Oldest forams	Phan	542	2, 200	(65)
Gonyaulacales	Gonyaulacaceae split	Phan	196	2, 10	(59)
Pennate diatoms	Oldest pennate	Phan	80	3, 5	(66)
Spirotrichs	Oldest tintinnids	Phan	444	2.5, 100	(67)
Trachaeophytes	Earliest trachaeophytes	Phan	425	4, 2.5	(68)
Vertebrates	<i>Haikouichthys</i>	Phan	520	3, 5	(69)
Animals	LOEMs, sponge biomarkers	Protero	632	2, 300	(70, 71)
Arcellinida	<i>Paleoarcella</i>	Protero	736	2, 300	(12)
Bilateria	<i>Kimberella</i>	Protero	555	2, 30	(72)
Chlorophytes	<i>Palaeastrum</i>	Protero	700	2.5, 300	(73)
Ciliates	Gammacerane	Protero	736	2.5, 300	(74)
Florideophyceae	Doushantuo red algae	Protero	550	2.5, 100	(75)
Red algae <sup>‡</sup>	<i>Bangiomorpha</i>	Protero	1174	3, 250	(11)

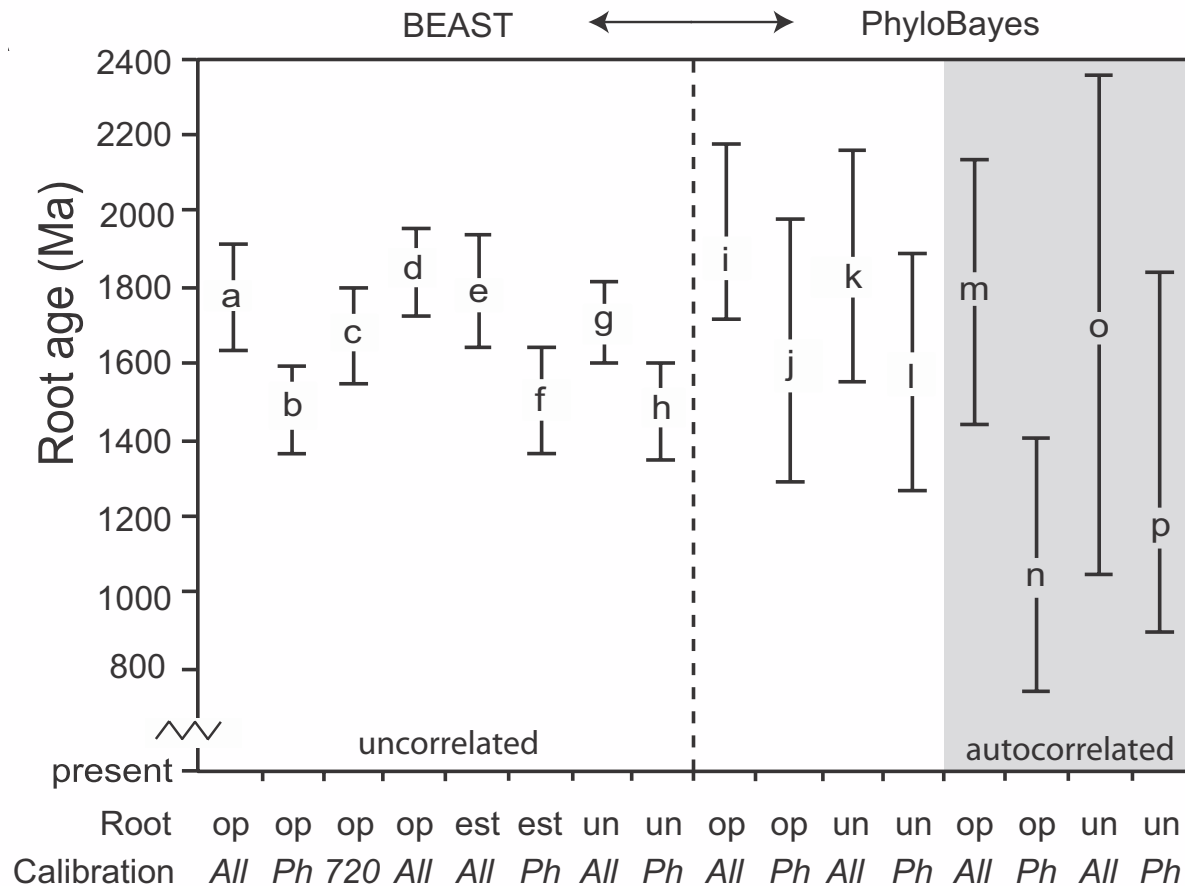
\*Eon: Phan, Phanerozoic; Protero, Proterozoic. Proterozoic calibrations are excluded from *Phan* analyses.

<sup>†</sup>Calibration constraints are specified for BEAST using a gamma distribution with a minimum date in Ma based on the fossil record parameters as indicated: min, minimum divergence data; dist, gamma prior distribution (shape, scale). See [Table S3](#) for details of PhyloBayes calibrations.

<sup>‡</sup>In the *All 720* analysis (c), the minimum age constraint for the red algae node is set to 720 Ma.

# Horloges moléculaires

impact des choix méthodologiques sur la datation

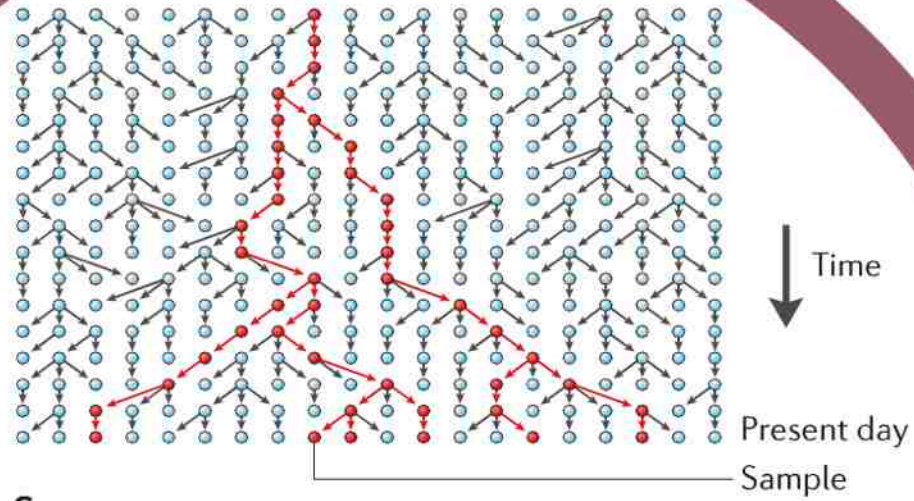


LECA 1000-1900 Ma

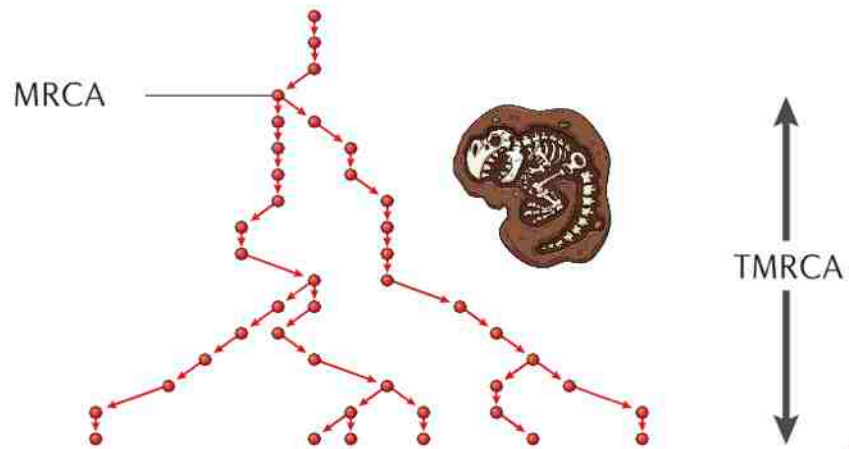


*Bangiomorpha*

Marjoram & Tavaré (2006)



**c**



fossiles  $\neq$  ancêtres