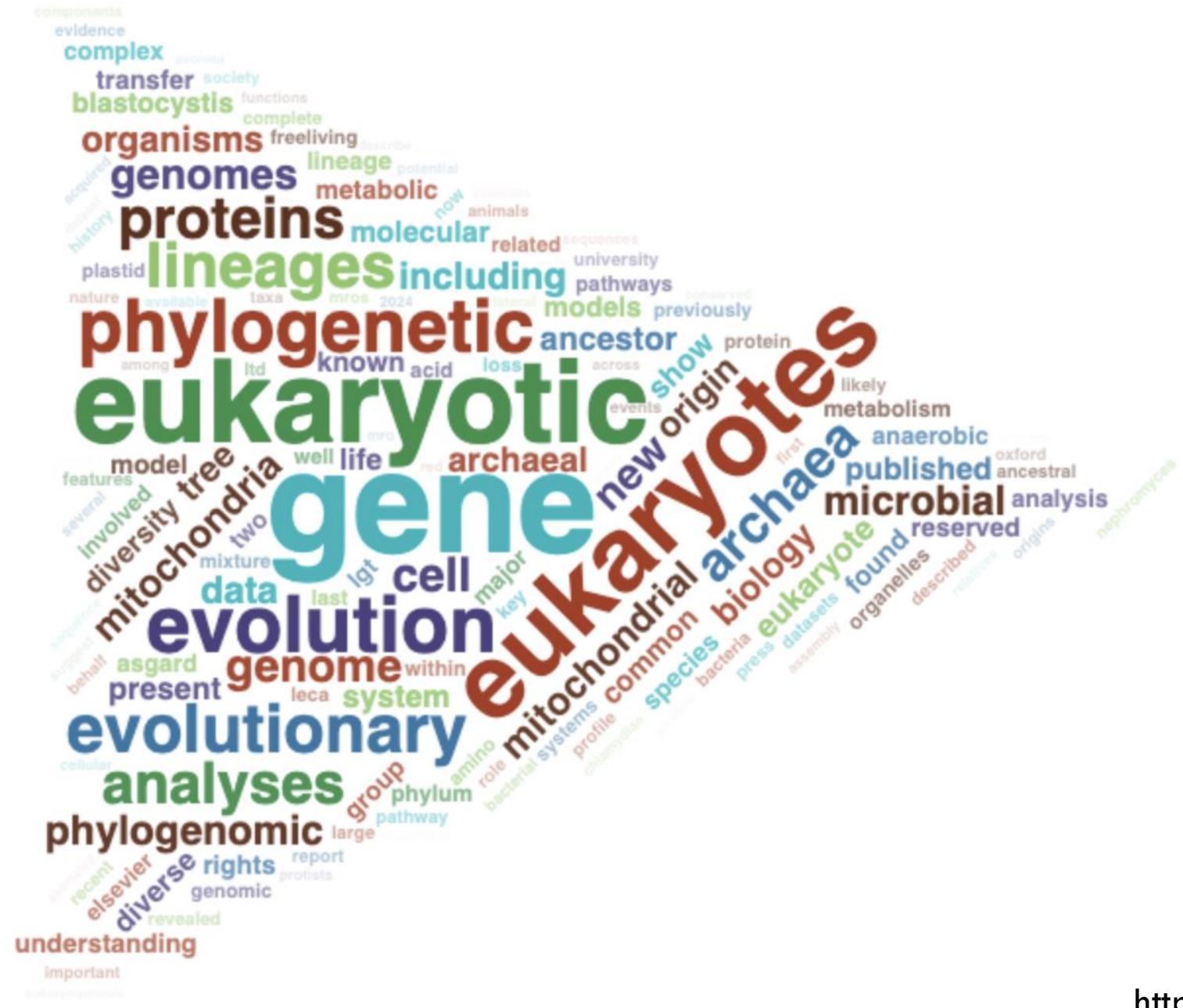


"Deep" phylogenetics

Laura Eme
Associate Professor
University of Rhode Island

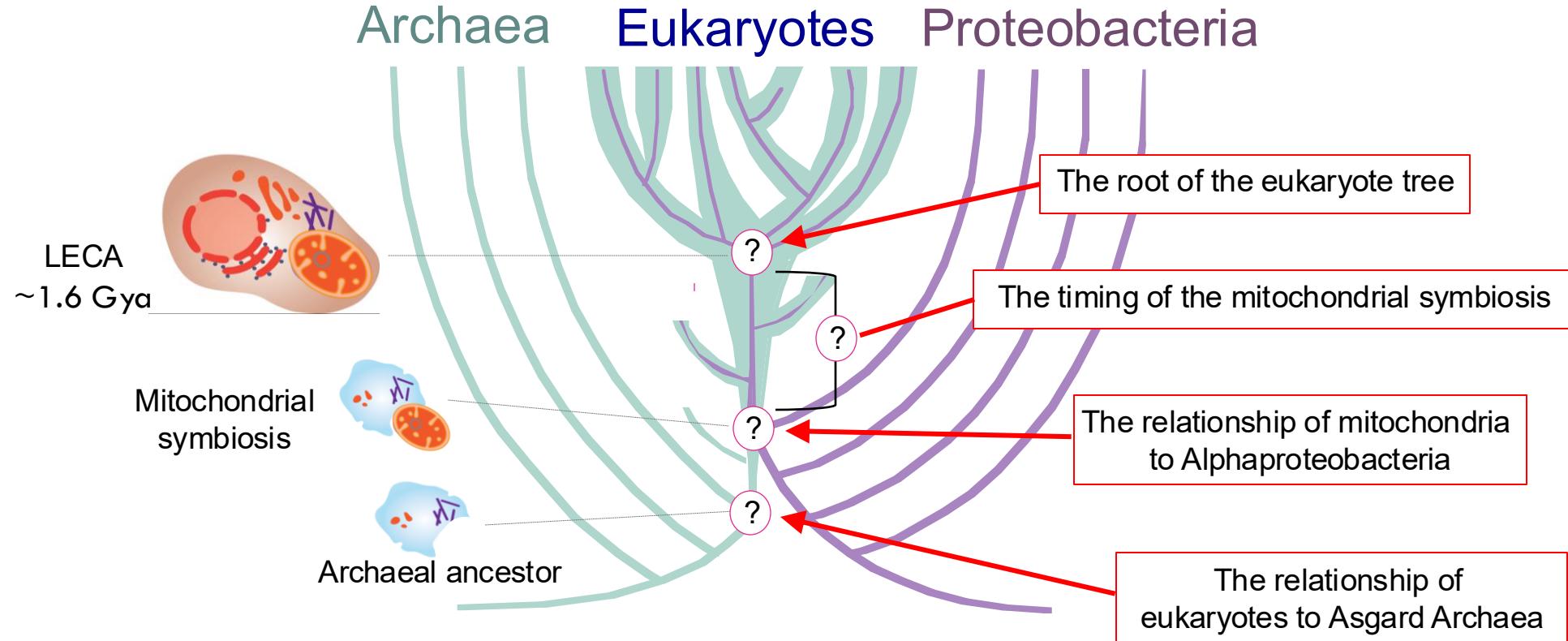
@lauraeme.bsky.social





<https://shiny.research.sfu.ca/>

Four major phylogenetic problems of eukaryogenesis



1 - Protein models of evolution

1.1 Empirical models

Code degeneracy

Glu-Gly-Ser-Ser-Trp-Leu-Leu-Leu-Gly-Ser

Glu-Gly-Ser-Ser-Tyr-Leu-Leu-Ile-Gly-Ser

Asp-Gly-Ser-Ala-Trp-Leu-Leu-Leu-Gly-Ser

Asp-Gly-Ser-Ala-Tyr-Leu-Leu-Ala-Gly-Ser

GAA-GGA-AGC-TCC-TGG-TTA-CTC-CTG-GGA-TCC

GAG-GGT-TCC-AGC-TAT-CTA-TTA-ATT-GGT-AGC

GAC-GGC-AGT-GCA-TGG-TTG-CTT-TTG-GGC-AGT

GAT-GGG-TCA-GCT-TAC-CTC-CTG-GCC-GGG-TCA

Protein sequence evolves slower than nucleotide

Code degeneracy

- Base composition bias can lead to large difference in codon usage
- Comparing protein sequences can reduce the compositional bias problem

Evolutionary models for amino acid changes

Typically

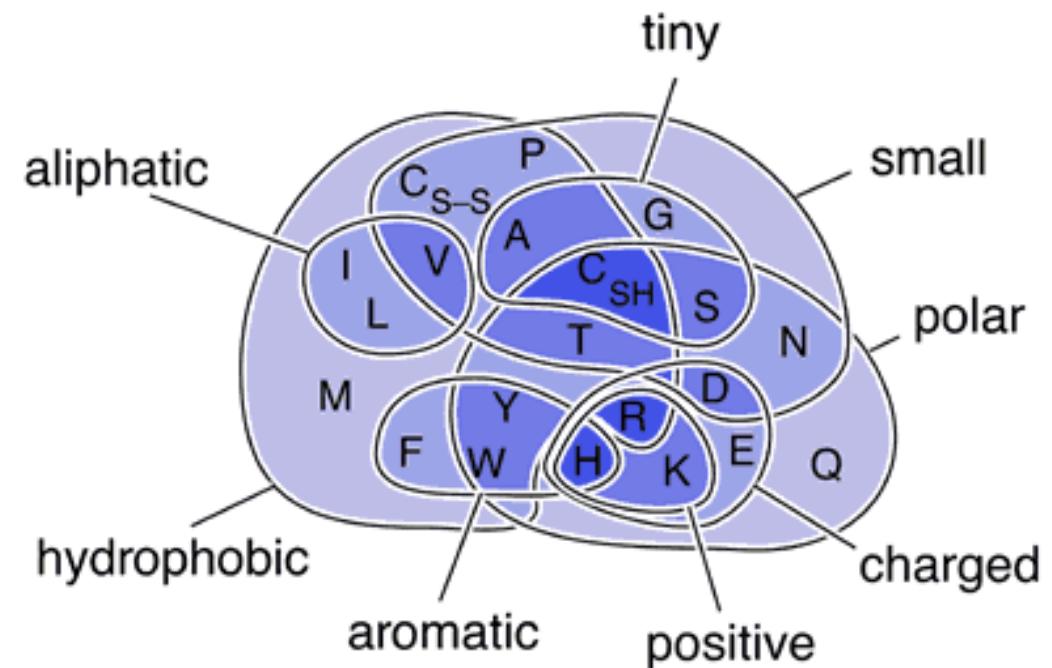
- A 20x20 rate matrix
- Assumes stationarity and reversibility

Amino acid physico-chemical properties

- AA can be categorized according to their physicochemical properties
- Major factor in protein folding (secondary, tertiary, quaternary structure)
- Key to protein functions (e.g., catalytic sites)

→ Major influence in pattern of amino acid mutations

Some amino acid changes are more commonly fixed than others



Empirical models: amino acid substitution matrices based on observed substitutions

Summarise the substitution patterns from a large number of existing alignments ('average' models)

Empirical models: amino acid substitution matrices based on observed substitutions

Summarise the substitution patterns from a large number of existing alignments ('average' models)



Raw data: observed changes in pairwise comparisons

seq.1 AIDESLIIASIATATI

| * | | * | | * | | * | | * | |

seq.2 AGDEALILASAATSTI

seq.1 A I D E S L I I A S I A T A T I

| * | | * | | * | | * | | * | |

seq.2 A G E E A L I L A S A A T S T I

	A	S	T	G	I	L	E	D
Raw matrix	A	3						
Symmetrical	S	2	1					
	T	0	0	1				
	G	0	0	0	0			
	I	1	0	0	1	2		
	L	0	0	0	0	1	1	
	E	0	0	0	0	0	0	1
	D	0	0	0	0	0	0	1

→ The larger the dataset, the better the estimates

Amino acid exchange matrices

$$\left(\begin{array}{ccccc} - & s_{1,2} & s_{1,3} & \dots & s_{1,20} \\ s_{1,2} & - & s_{2,3} & \dots & s_{2,20} \\ s_{1,3} & s_{2,3} & - & \dots & s_{3,20} \\ \dots & \dots & \dots & \dots & \dots \\ s_{1,20} & s_{2,20} & s_{3,20} & \dots & - \end{array} \right)$$

$$X \operatorname{diag}(\pi_1, \dots, \pi_{20}) = Q \text{ matrix}$$

Q Rate matrix

s_{ij} Exchangeabilities of amino acid pairs ij

$s_{ij} = s_{ji}$ Time reversibility (usually)

π_i Stationarity of amino acid frequencies
(typically the observed proportion of residues in the dataset)

Empirical models

- Summarise the substitution patterns from a large number of existing alignments ('average' models)
- Different substitution matrices come from:
 - Selection of specific proteins
 - Globular proteins vs membrane proteins
 - Mitochondrial proteins, viral proteins...
 - Range of sequence similarities used
 - Counting methods
 - On a tree
 - Pairwise comparison from an alignment

Empirical models

Dayhoff (Dayhoff et al., 1978): Nuclear encoded genes (~100 proteins) → PAM matrices

JTT (Jones et al., 1992): 59,190 point mutations from 16,300 proteins from membrane spanning segments

Limitation: for less similar sequences, no linearity between observed and real substitution rate (hidden substitutions)

Empirical models

Dayhoff (Dayhoff et al., 1978): Nuclear encoded genes, ~100 proteins → PAM matrices

JTT (Jones et al., 1992): 59,190 point mutations from 16,300 proteins from membrane spanning segments

WAG (Whelan and Goldman, 2001): General matrix

LG (Le and Gascuel, 2008): General matrix

The WAG matrix (2001)

- Globular protein sequences
 - 3,905 sequences from 182 protein families
- Produced a phylogenetic trees for every family and used maximum likelihood to estimate the relative rate values in the rate matrix (i.e., maximizes the overall lnL over 182 different trees)
- Better fit of the model with most data (significant improvement of the tree lnL when compared to PAM or JTT matrices)
- Can be used for (more) distant homologues

Further improvements: the LG matrix (2008)

- Used the same phylogenetic approach as WAG
- Further refine the method by adding the variability of evolutionary rates across sites when estimating the matrix and increase the number of sequences used
- Better fit of the model with most data (significant improvement of the tree lnL when compared to WAG and other matrices)

Empirical models

Dayhoff (Dayhoff et al., 1978): Nuclear encoded genes, ~100 proteins → PAM matrices

JTT (Jones et al., 1992): 59,190 point mutations from 16,300 proteins from membrane spanning segments

WAG (Whelan and Goldman, 2001): General matrix

LG (Le and Gascuel, 2008): General matrix

Mtrev24 (Adachi and Hasegawa, 1996) : Mitochondrial (vertebrates)

Mtmam (Yang et al., 1998): Mitochondrial (mammals)

mtART (Abascal et al., 2007): Mitochondrial (Arthropoda)

CpRev (Adachi et al., 2000): Chloroplast

VT (Müller and Vingron, 2000): General matrix

RtRev (Dimmic et al., 2002): Retrovirus

DayhoffDCMUT (Kosiol and Goldman, 2005): Revised Dayhoff matrix

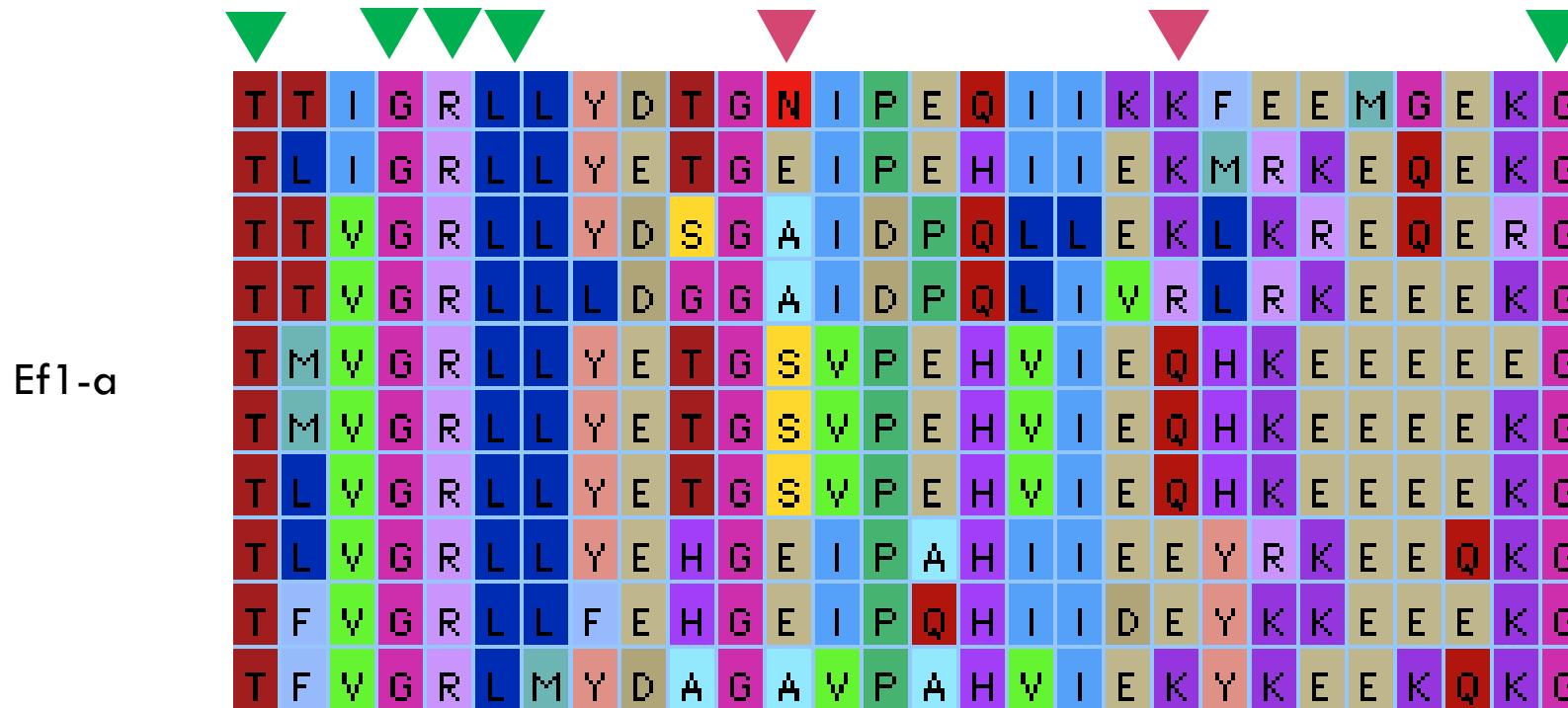
(and more...)

Summary

- Many amino acid rate matrices exist
- One should make a rational choice (as much as possible):
 - How was the rate matrix produced?
 - What are the structural features of the sequences that you are analyzing? Globular/membrane protein? Overall level of sequence identity of the compared sequences? Specific compositional bias (mitochondrial proteins matrix: mtREV24; Transmembrane domains: PHAT)?
 - ModelTest, ModelFinder (IQtree), ProtTest... to compare models

Rate heterogeneity parameter

- Not all sites “evolve” at the same speed depending on how it impacts function



Rate heterogeneity parameter

- Discreticized Gamma distribution (+G)
 - Default is usually 4 categories but can be set to be more (but more computationally intensive)

Rate heterogeneity parameter

- Discreticized Gamma distribution (+G)
 - Default is usually 4 categories but can be set to be more (but more computationally intensive)
- FreeRate model (+R)
 - Does not follow a parametric distribution
 - Not all categories will have the same number of sites
 - More realistic but more computationally intensive
 - Typically fits data better than the +G model and is recommended for analysis of large data sets

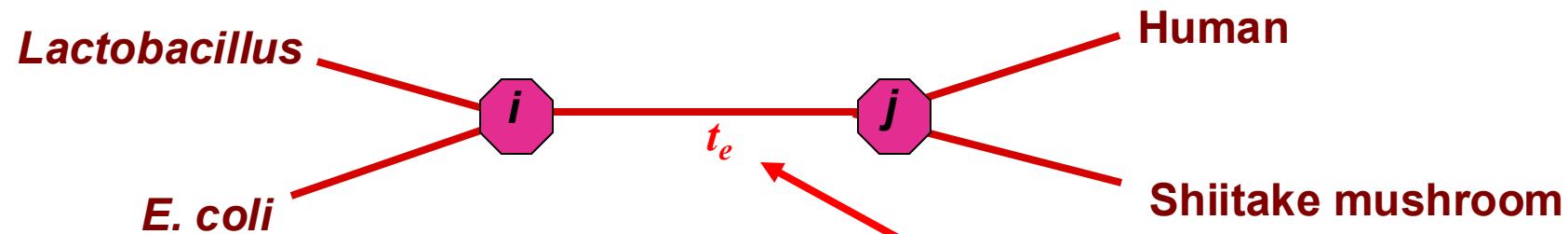
1.2 Fully parameterized time-reversible model

GTR (General time reversible)

- One can generate a dataset-specific model
 - All parameters of the Q matrix are estimated from your data (exchangeabilities and equilibrium frequencies)
 - GTR20: General time reversible model for amino-acids: 189 rate parameters!
- *WARNING* Parameter-rich: parameter estimates might not be reliable if made on short alignments (not enough information)

1.3 Mixture models

Your model is giving you the probability of going from amino acid i to j at site x , evolving at rate r_v on branch t_e

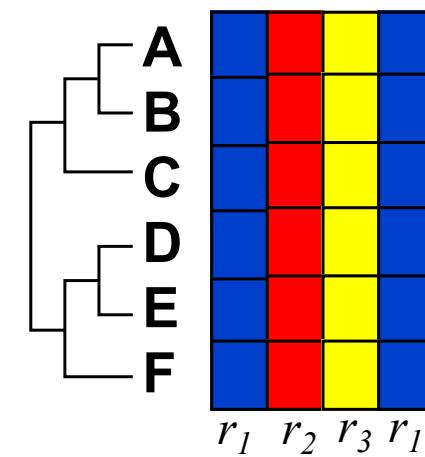
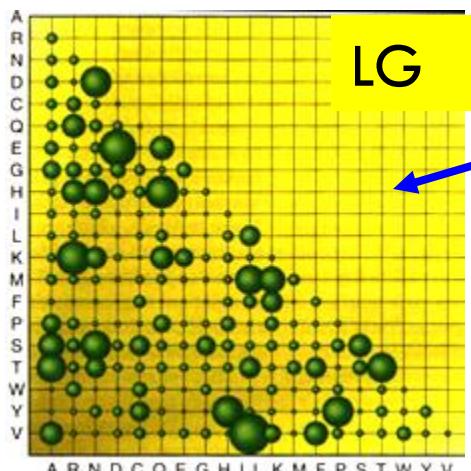


$$P(j | i; t) = [\exp(Q \square t_e \square r_v)]_{ij}$$

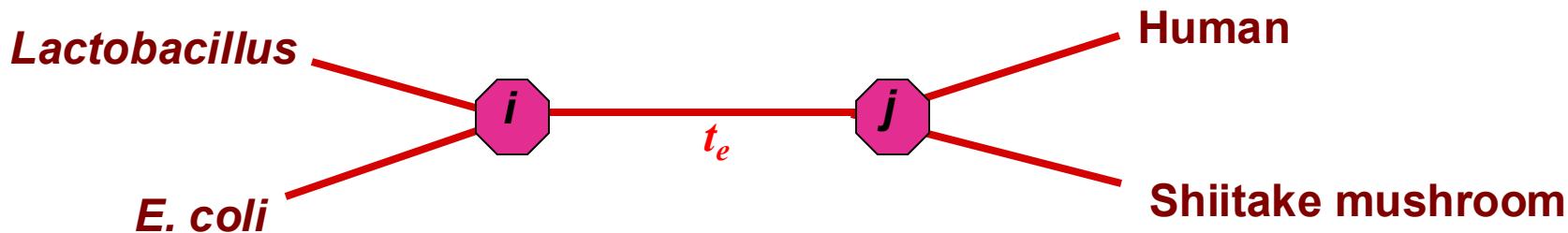
For $i \neq j$:

$$q_{ij} = r_{ij} \square \pi_j$$

$$\Pi = \begin{bmatrix} \square \pi_A & 0 & 0 & 0 \\ \square 0 & \pi_R & 0 & 0 \\ \square 0 & 0 & \dots & 0 \\ \square 0 & 0 & 0 & \pi_v \end{bmatrix}$$



Your model is giving you the probability of going from amino acid i to j at site x , evolving at rate r_v on branch t_e

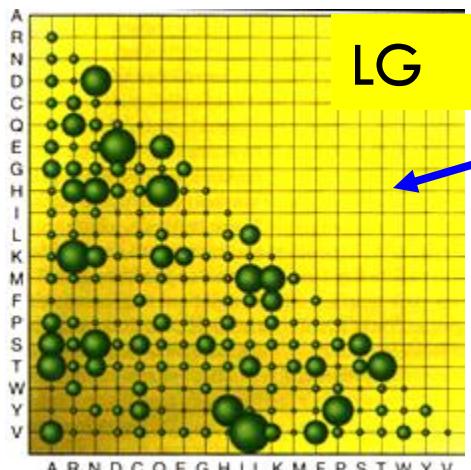


$$P(j | i; t) = [\exp(Q \square t_e \square r_v)]_{ij}$$

For $i \neq j$:

$$q_{ij} = r_{ij} \square \pi_j$$

$$\Pi = \begin{bmatrix} \square \pi_A & 0 & 0 & 0 \\ \square 0 & \pi_R & 0 & 0 \\ \square 0 & 0 & \dots & 0 \\ \square 0 & 0 & 0 & \pi_v \end{bmatrix}$$



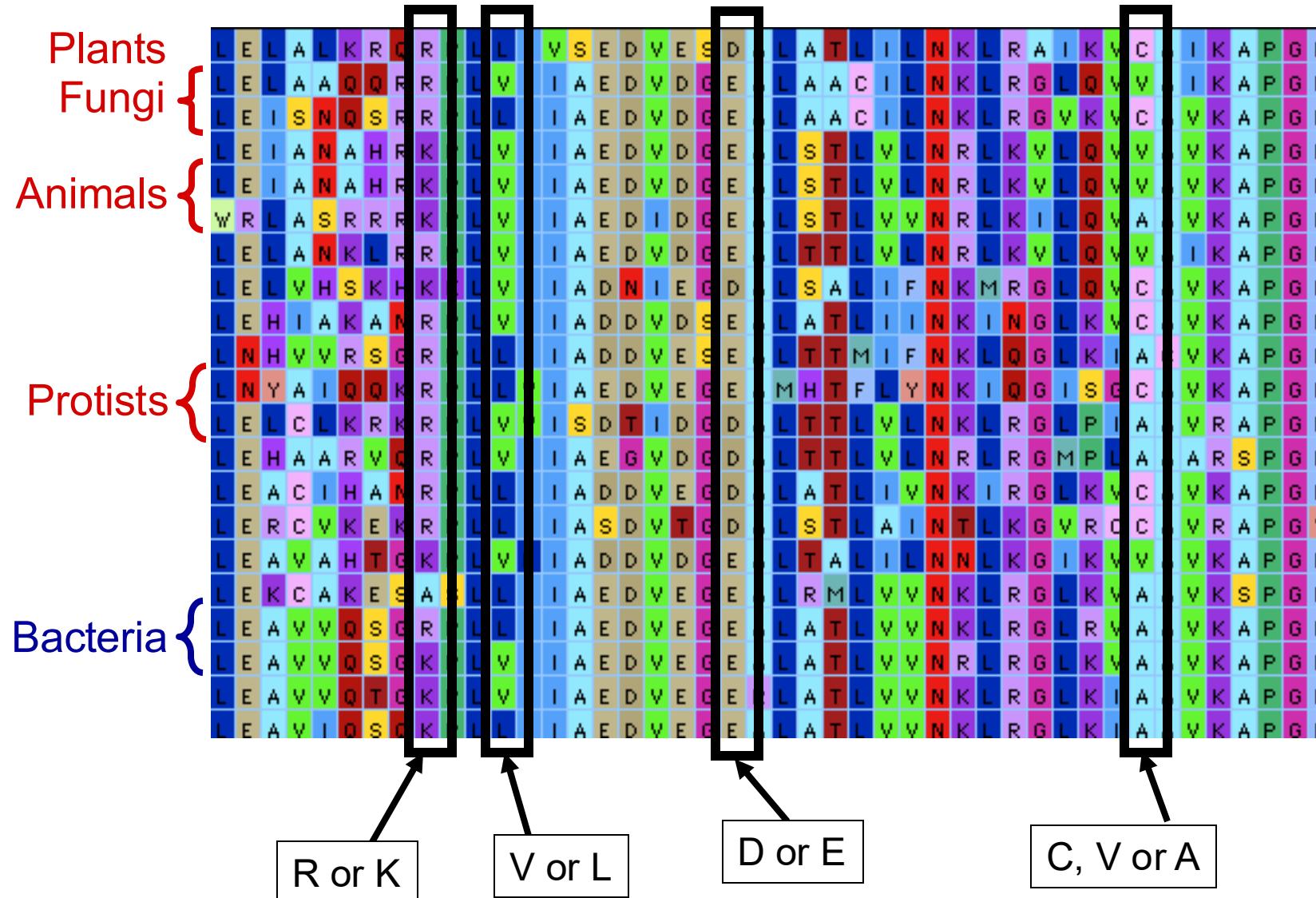
Assumptions

- different sites in protein and organisms all evolve according to the same general ‘rules’
- i.e. rate matrices (R ’s) and frequencies (Π ’s) are the same for all sites and branches

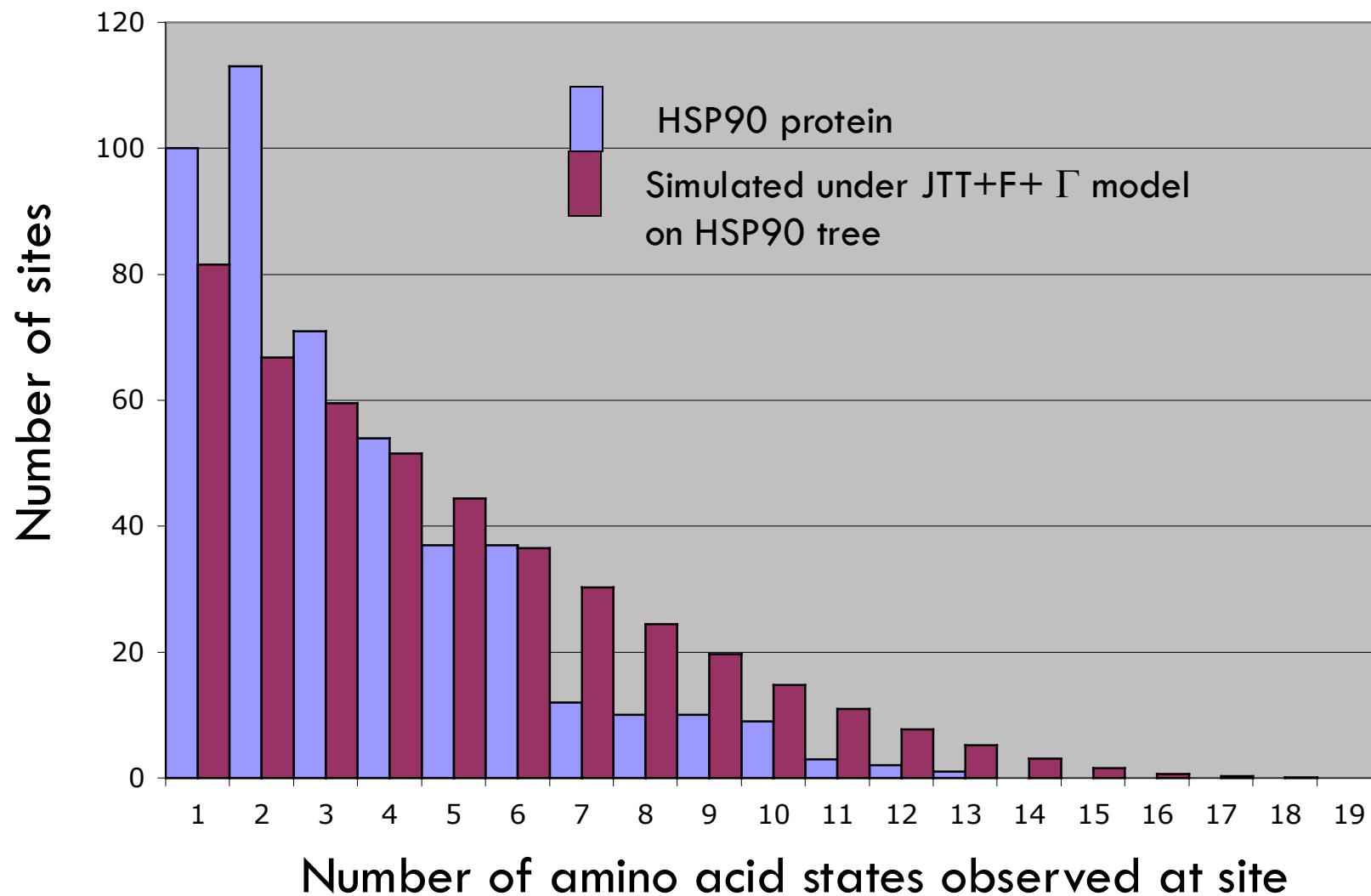
The problem...

- Such models are a dramatic over-simplification of what is really going on
 - Average over sites, average over different organisms, average across protein families
- Sites in proteins can change function over time
 - sites under negative selection \leftrightarrow neutral \leftrightarrow positive selection
- Every amino acid site in a protein has a unique structural/functional context
 - Hydrophobicity, polarity, charge, size, functional group, etc.
 - Different sites have different exchangeabilities
 - Different frequencies of AAs occur at different sites

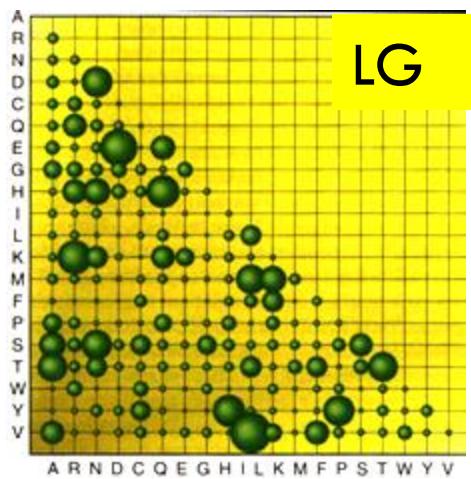
Evolution of chaperonin 60 over ~1.5 billion years



Distribution of the number of different amino acids at aligned sites

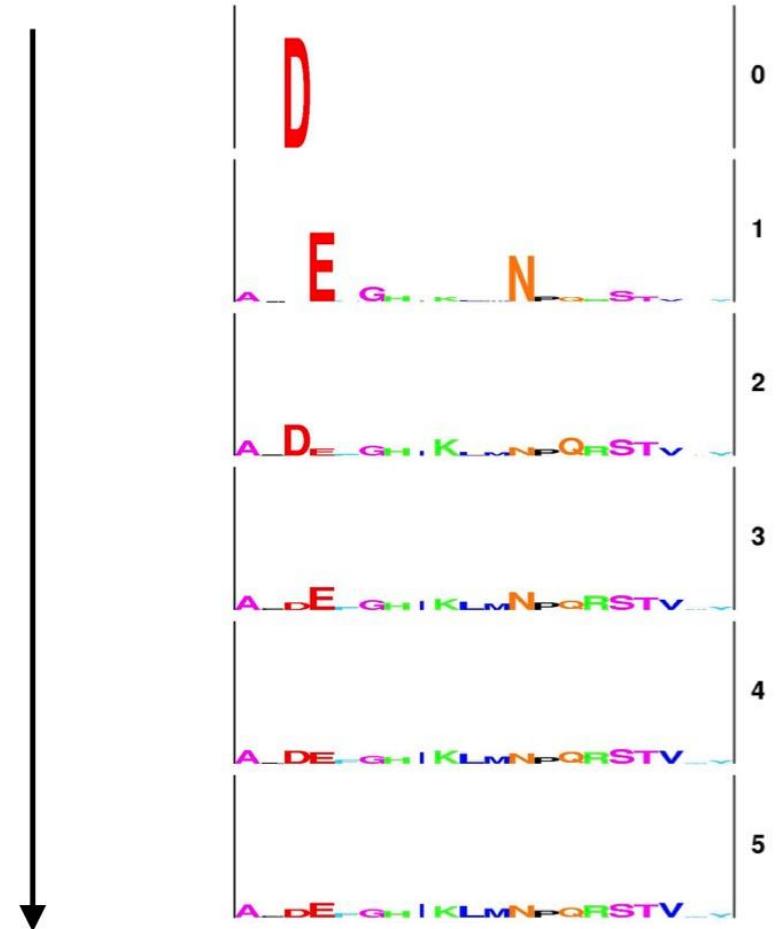


Starting at a D with a site homogeneous matrix (LG+F)

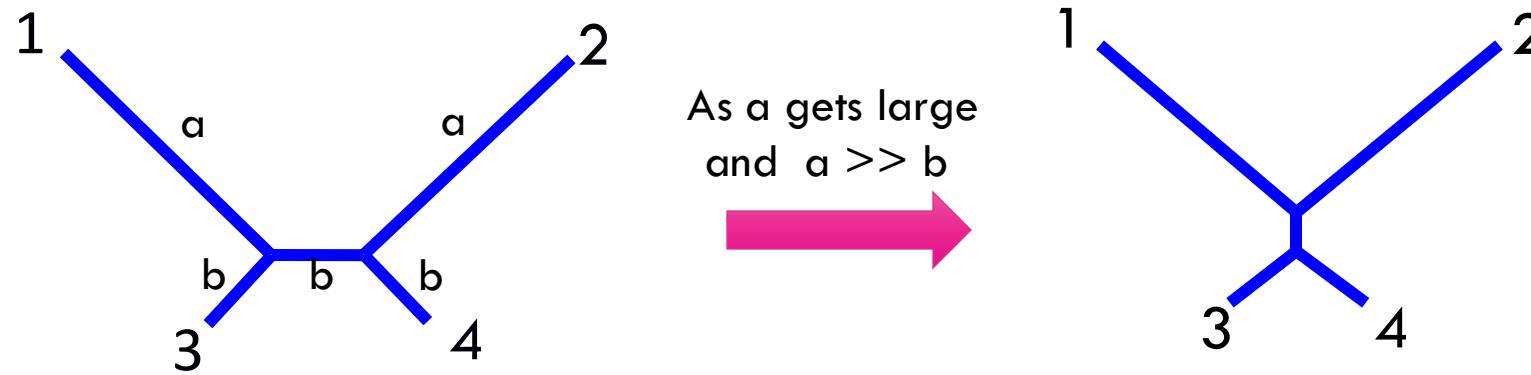


$$\Pi = \begin{bmatrix} \pi_A & 0 & 0 & 0 \\ 0 & \pi_R & 0 & 0 \\ 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & \pi_v \end{bmatrix}$$

substitutions



What happens to phylogenetic estimation when you ignore site-heterogeneity?



Long branch attraction

Susko et al. (2004) *Mol. Biol. Evol.*

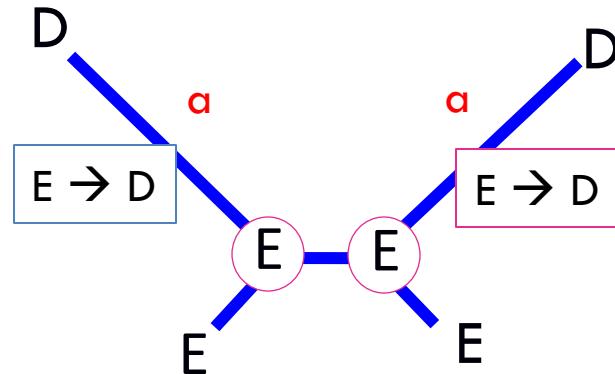
Lartillot and Philippe (2007) *BMC Evol. Biol.*

Wang et al. (2008) *BMC Evol. Biol.*

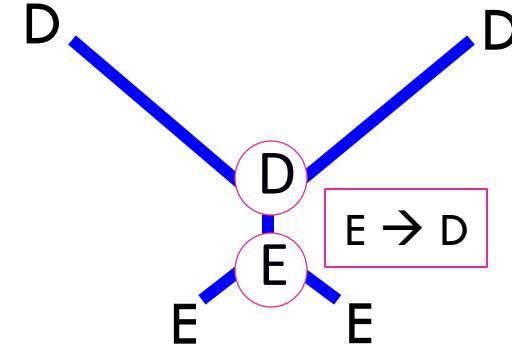
Roger and Susko (2021) *Systematic Bio*

Why long branch attraction (LBA)?

TRUE TREE



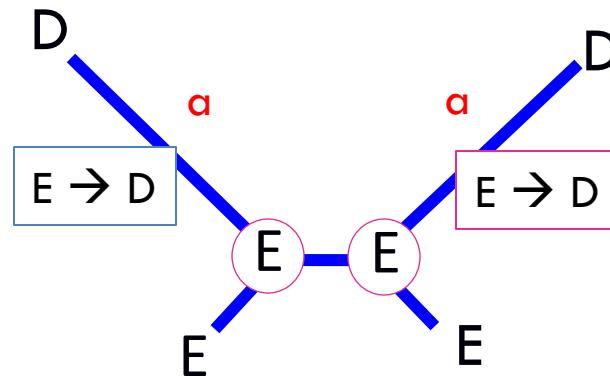
LBA TREE



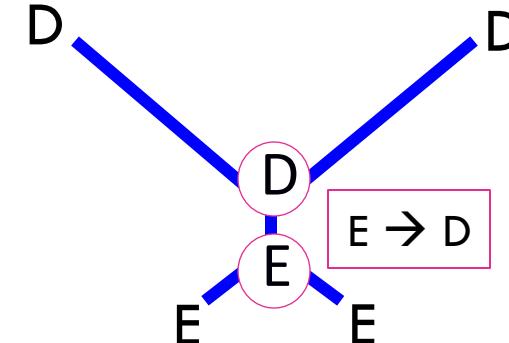
Under site homogeneous model (LG), the probability of converging on the same state ($E \rightarrow D$) twice is pretty low:
→ if branch-length a is really long, then $P(\text{convergence})_{\text{LG}} \approx \pi_D^2 = (0.057)^2 = 0.0032$

Why long branch attraction (LBA)?

TRUE TREE



LBA TREE



Under site homogeneous model (LG), the probability of converging on the same state ($E \rightarrow D$) twice is pretty low:

→ if branch-length a is really long, then $P(\text{convergence})_{\text{LG}} \approx \pi_D^2 = (0.057)^2 = 0.0032$

Under a site-specific model where you can only be D or E (with equal frequency of 0.5):

→ $P(\text{convergence})_{\text{ss}} \approx \pi_D^2 = (0.5)^2 = 0.25$

Mixture models

- Standard protein substitution models: single Q matrix
- Mixture models: combine several amino-acid replacement matrices
- Same principle as rate heterogeneity mixture models
 - For each site, its likelihood is the **sum of its weighted likelihood under each Q matrix** that are part of the mixture model

$$L_i = p(\mathbf{y}_i | r_1)p(r_1) + p(\mathbf{y}_i | r_2)p(r_2) + \cdots + p(\mathbf{y}_i | r_k)p(r_k)$$

Weight of the
rate class

Mixture models: terminology warning

- Different kinds of mixture models!
- Rate-category mixture model (see Dave's lecture)
- **Usually people refer to mixture of amino-acid replacement matrices**
- Mixtures can be apply to any part of the model (e.g., branch lengths)

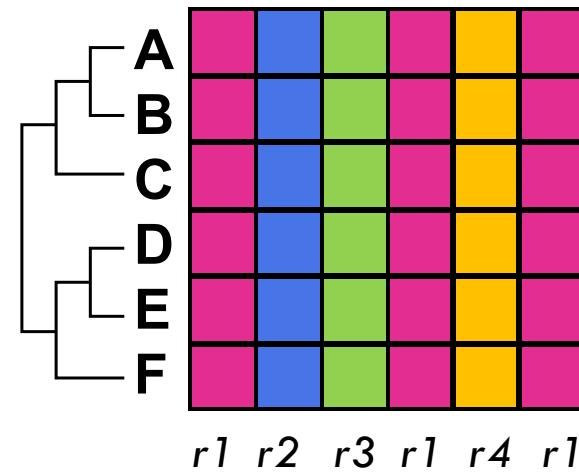
LG4M and LG4X mixture models

“the variability of evolutionary rates corresponds to one of the most apparent heterogeneity factors among sites, and **there is no reason to assume that the substitution patterns remain identical regardless of the evolutionary rate**” Le, Dang, Gascuel 2012

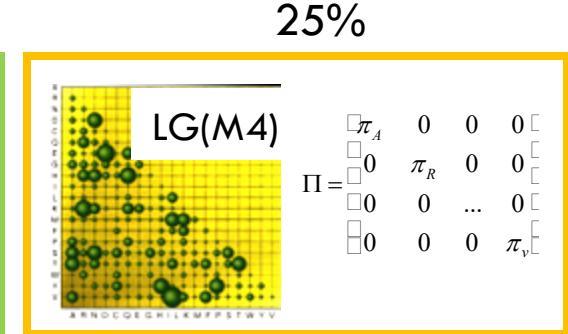
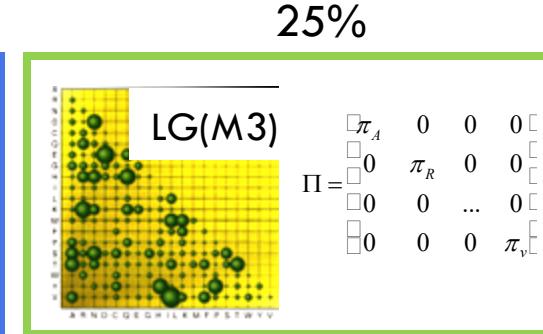
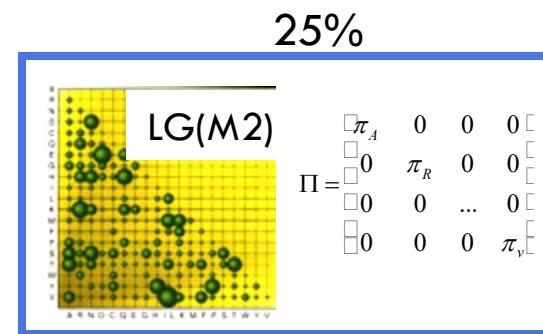
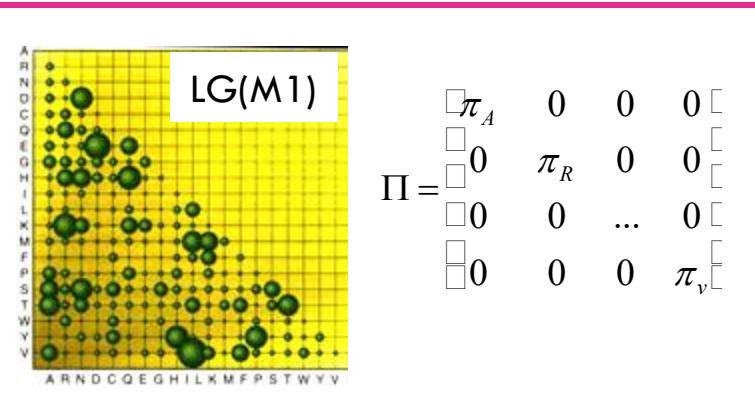
Standard LG+G model:
only the global rate differs from one category to another

LG4M and LG4X mixture models

LG4M: each gamma rate category gets its own Q matrix (i.e., each of the 4 gamma-distributed rate category gets its own amino acid equilibrium distributions and exchangeabilities)



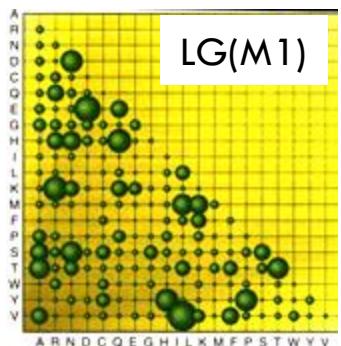
25%



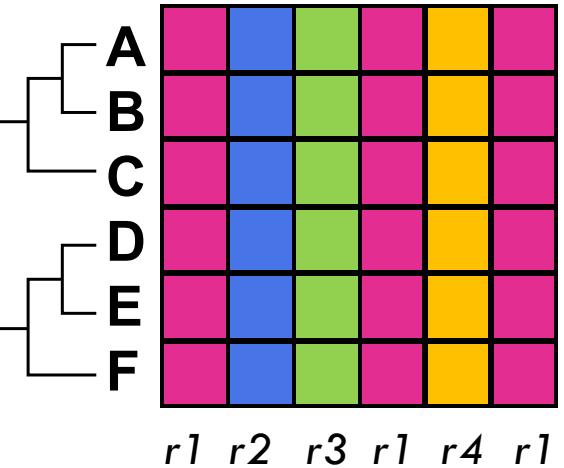
LG4M and LG4X

LG4X: each rate category gets its own Q matrix BUT rates and weights are left out of the gamma distribution assumption

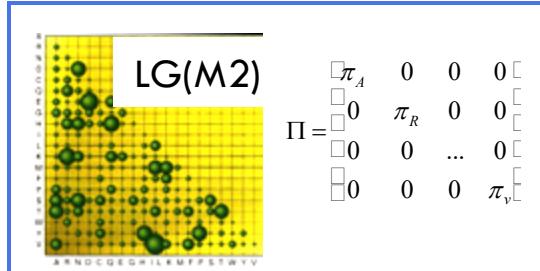
50%



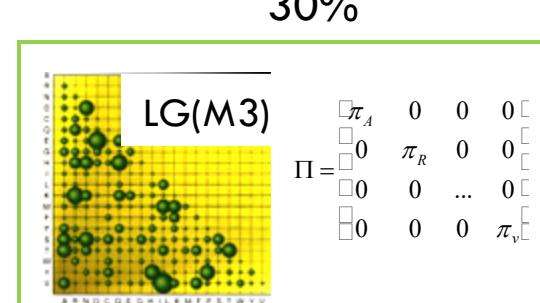
$$\Pi = \begin{bmatrix} \pi_A & 0 & 0 & 0 \\ 0 & \pi_R & 0 & 0 \\ 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & \pi_v \end{bmatrix}$$



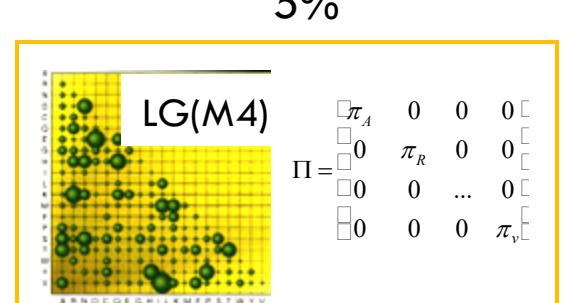
15%



15%

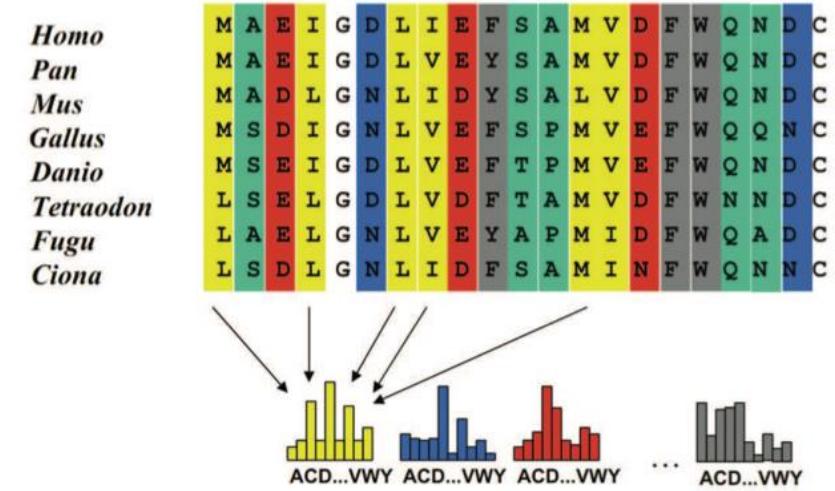


30%



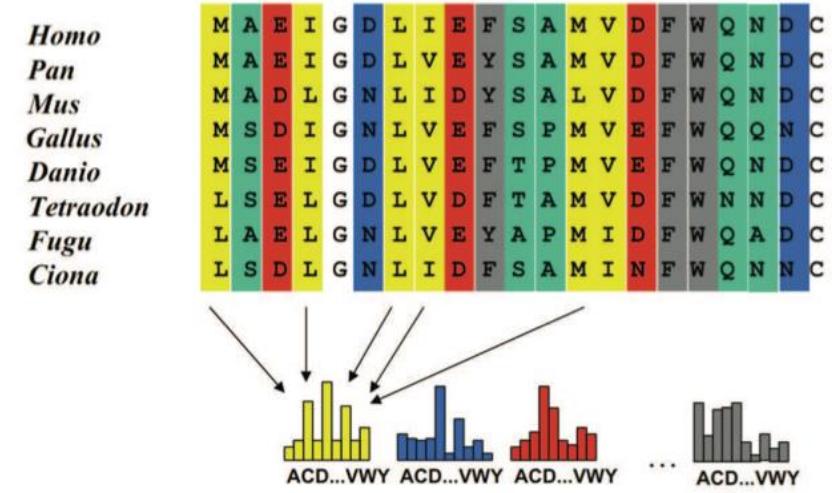
5%

The CAT model



- Bayesian framework only
- Free number of profiles in the mixture model (estimated during the Bayesian procedure). “Infinite mixture model”
- Each profile corresponds in practice to a *biochemical profile*: only a small number of AA are highly probable, while the frequency of all others will be ~ 0 .

The CAT model



- CAT-Poisson: very simple amino-acid replacement process (R matrix).
Each time a substitution event occurs, a new amino-acid is chosen at random, according to the probabilities defined by the profile (Poisson or proportional amino-acid replacement process).
Eg., any AA has the same probability to mutate to a Valine.
- CAT-GTR: GTR exchangeability matrix with 189 parameters!

C10, C20, ..., C60 mixture models

- 10, 20, 30, 40, 50, 60-profile mixture models are approximations of the CAT model for ML
- 10 (20, 30...) different pre-computed (empirical) Q matrices that correspond to 10 (20, 30...) most-common types of biochemical profiles in proteins
- By default, assume Poisson AA replacement but can be combined with empirically estimated exchangeabilities, such as from the LG matrix.
For example: LG+C10

$$q_{ij} = r_{ij} \square \pi_j$$

Problem with mixture models

- As the number of sites and proteins increases the computational cost becomes prohibitive
 - For an ML analysis of 104 taxa and ~90,000 sites (350 proteins concatenated) LG+C60+F+G model takes >350 GB of RAM and ~3 weeks on 12 cores to estimate the **ML tree** using IQTREE v. 1.5
 - **5.5 years to do true bootstrap analysis**
- PMSF (Posterior Mean Site Frequency) approximation: transforming a mixture model into a ‘simple’ model (Wang, Bui, Susko, Roger *Systematic Biology* 2018)

PMSF (Posterior Mean Site Frequency) model

Implemented in IQtree

- 1) Reconstruct an ML tree under a “good” model = guide tree
- 2) Using the guide tree, estimate, **for each site x**, the posterior probability of each amino-acid class c (e.g.: C1, C2, ..., C60)

Posterior probability of ‘class c’ at site x

$$P(c|x) = \frac{w_c \times P(x|c)}{\sum_c w_c \times P(x|c)}$$

- 3) For each site x, estimate the **posterior mean frequency of each amino acid j**

Posterior mean frequency of amino acid j at site x over all c classes ($f_{j,x}$)

$$f_{j,x} = \sum_c f_{j,c} \times P(c|x)$$

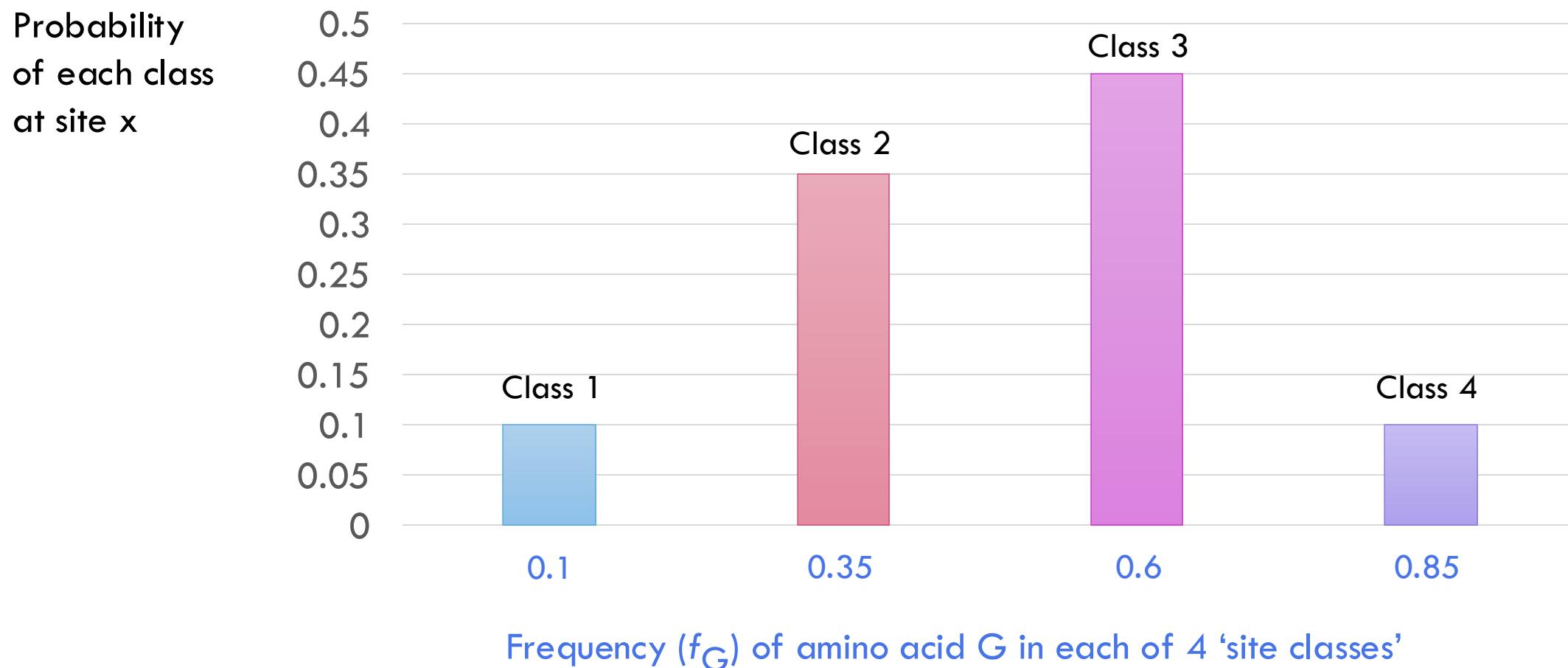
Sum over all classes

Freq of AA j for class c Prob of class c at site x

Example: Posterior mean site frequency for 'G' at a given site x, with a 4 class mixture model

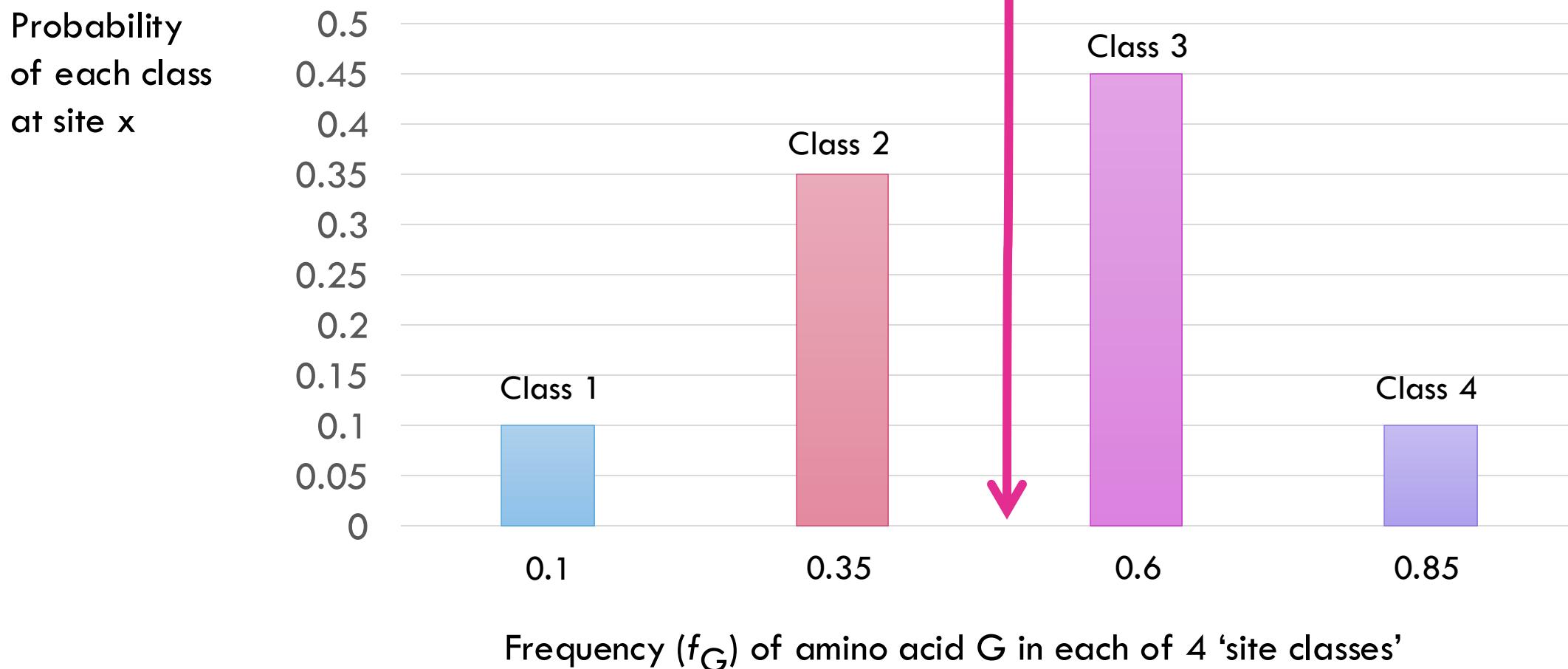


Example: Posterior mean site frequency for 'G' at a given site x, with a 4 class mixture model



E.g.: Posterior mean site frequency for 'G' at a given site x, with a 4 class mixture model

$$E[f_G] = (0.1 \times 0.1) + (0.35 \times 0.35) + (0.6 \times 0.45) + (0.85 \times 0.1) = 0.5$$



PMSF (Posterior Mean Site Frequency) model

- 1) Reconstruct an ML tree under a ‘reasonably good’ model
- 2) Using the ML tree, estimate, for each site x , the posterior probability of each amino-acid class c of your preferred mixture model (e.g.: C60)

Posterior probability of ‘class c ’ at site x

$$P(c|x) = \frac{w_c \times P(x|c)}{\sum_c w_c \times P(x|c)}$$

- 3) For each site x , estimate the posterior mean frequency of each amino acid j

Posterior mean frequency of amino acid j
at site x over all c classes ($f_{j,x}$)

$$f_{j,x} = \sum_c f_{j,c} \times P(c|x)$$

- 4) Now, every site x has its own $\Pi =$

$$\begin{bmatrix} \pi_A & 0 & 0 & 0 \\ 0 & \pi_R & 0 & 0 \\ 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & \pi_v \end{bmatrix}$$

PMSF (Posterior Mean Site Frequency) model

5) You estimate the ML tree using these pre-computed site-specific Q matrices: LG exchangeabilities + custom frequencies

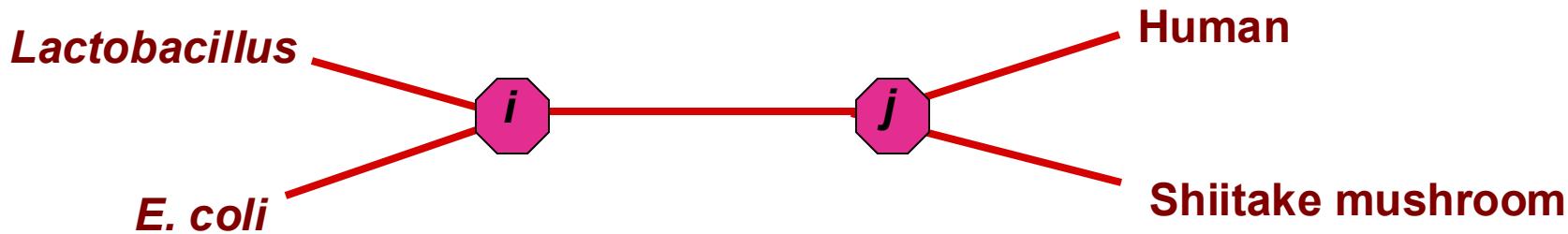
- Equivalent to LG+F, where F would be different for every site
- Barely more computationally intensive than using the ‘native’ LG matrix
- Bootstrapping is dramatically faster

Take home (for this part)

- Models are idealizations of the actual process of protein evolution
- Model misspecification (e.g. single-matrix models) often means systematic error (LBA)
- Mixture models deal with site-specific heterogeneity but are computationally expensive
- PMSF models provide a viable alternative for bootstrap analyses

Other types of mixture models

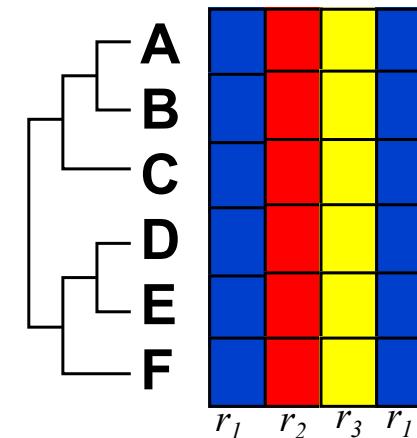
Probability of going from amino acid i to j
at site x , evolving at rate r_v on branch t_e



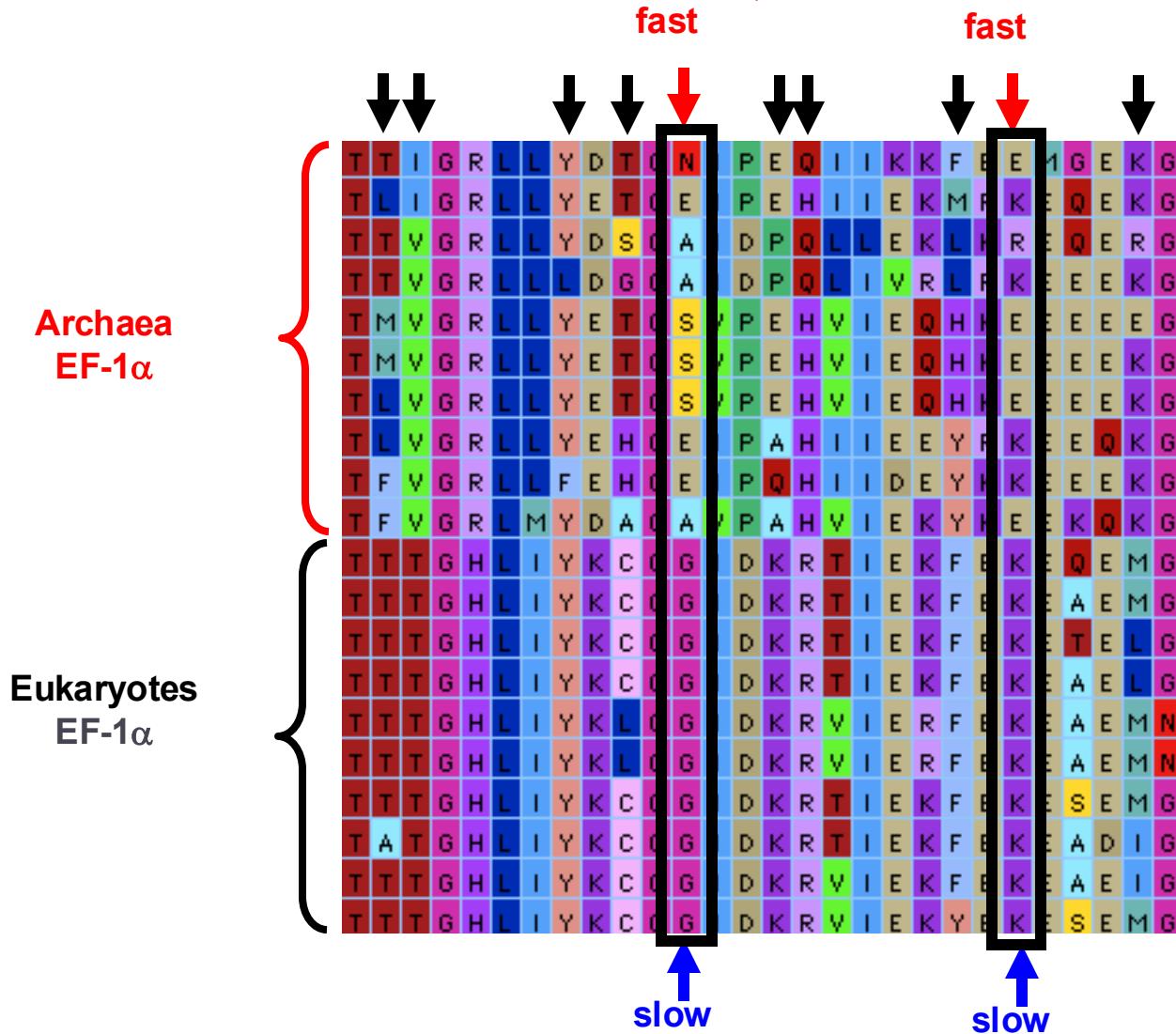
$$P(j | i; t) = [\exp(R \square \prod \square t_e \square r_v)]_{ij}$$

Assumptions

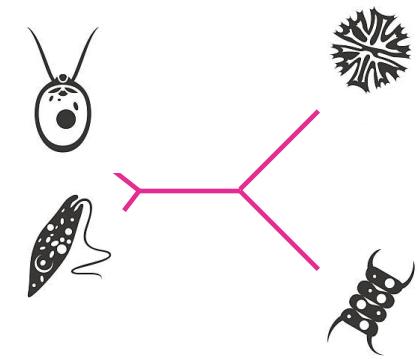
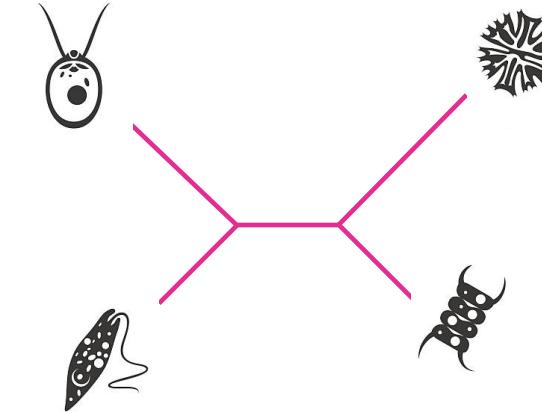
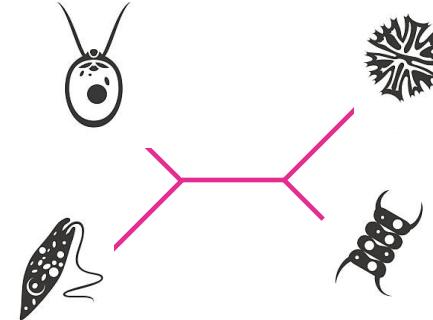
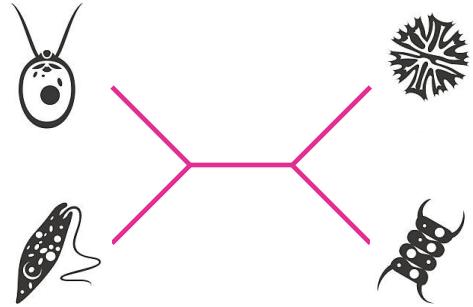
- ‘fast-evolving’ positions are always fast and slow-evolving positions are always slow
- Sites have the same rate of evolution (r_v) on different branches of tree



Changing rates of evolution at sites in different parts of the tree of life (=heterotachy)



Changing rates of evolution at sites in different parts of the tree of life (=heterotachy)



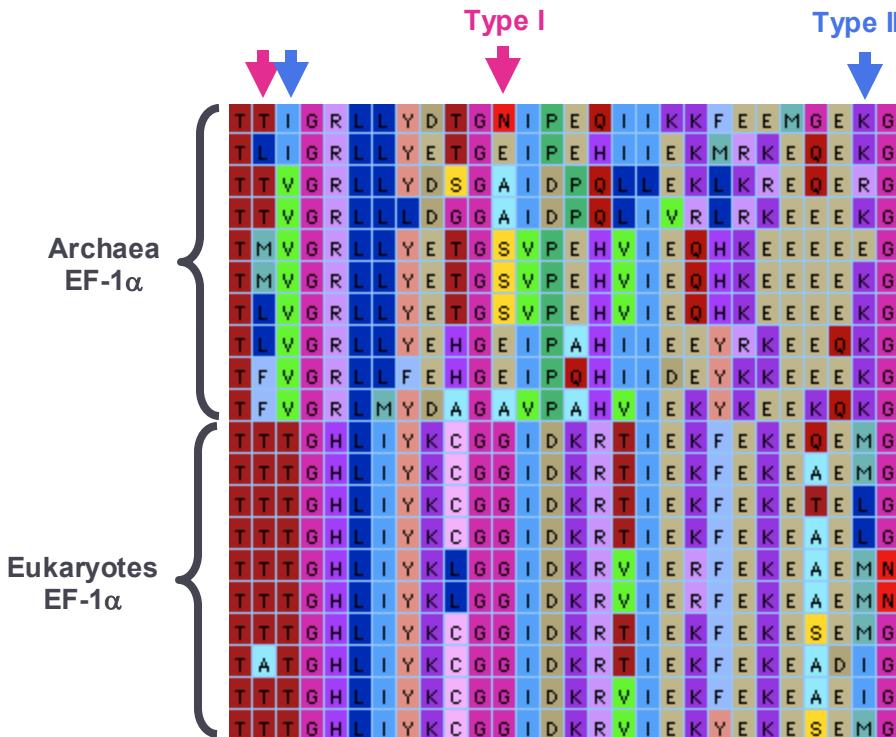
Models that deal with heterotachy (changing site rates across the tree)

- Covarion models (cf Joe's lecture)
 - Allow the sites “switch” between high rates and low rates over the tree
 - Computationally intensive
- Rate-shift models
 - Allows rates at many different sites to change abruptly on one branch
- Mixture of branch-length models
 - Allows different branch-lengths for different sites (e.g. GHOST model in IQtree)

Functionally divergent sites generate heterotachy

Functional shifts (functional divergence)

- Type I: 'rate-shifting' sites (sites that are conserved in one phylogenetic sub-group but not another).
- Type II: 'conserved-but-different' (conservation within both sub-groups of a phylogenetic tree but for amino acids with differing physico-chemical properties).



Functionally divergent sites generate heterotachy

Functional shifts (functional divergence)

- Type I: 'rate-shifting' sites (sites that are ~~discovered in one~~ FD sites violate homogeneity assumption phylogenetic sub-group but not another).
- Type II: conserved-but-different' (conservation within both sub-groups of a phylogenetic tree but for amino acids with differing physico-chemical properties).



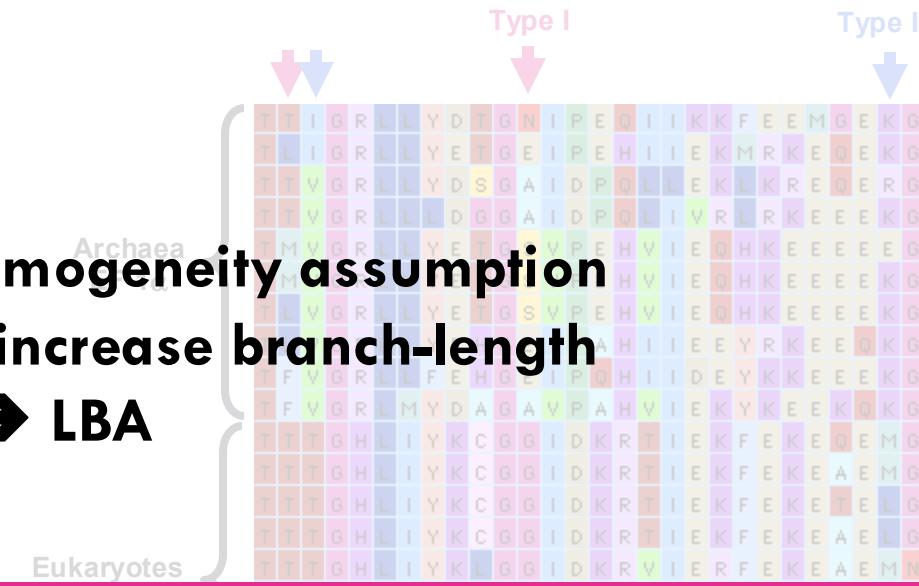
Functionally divergent sites generate heterotachy

Functional shifts (functional divergence)

- Type I: 'rate-shifting' sites (sites that are conserved in one phylogenetic sub-group but not another)
- Type II: conserved-but-different' (conservation within both sub-groups of a phylogenetic tree but for

FD sites violate homogeneity assumption and artefactually increase branch-length

→ LBA

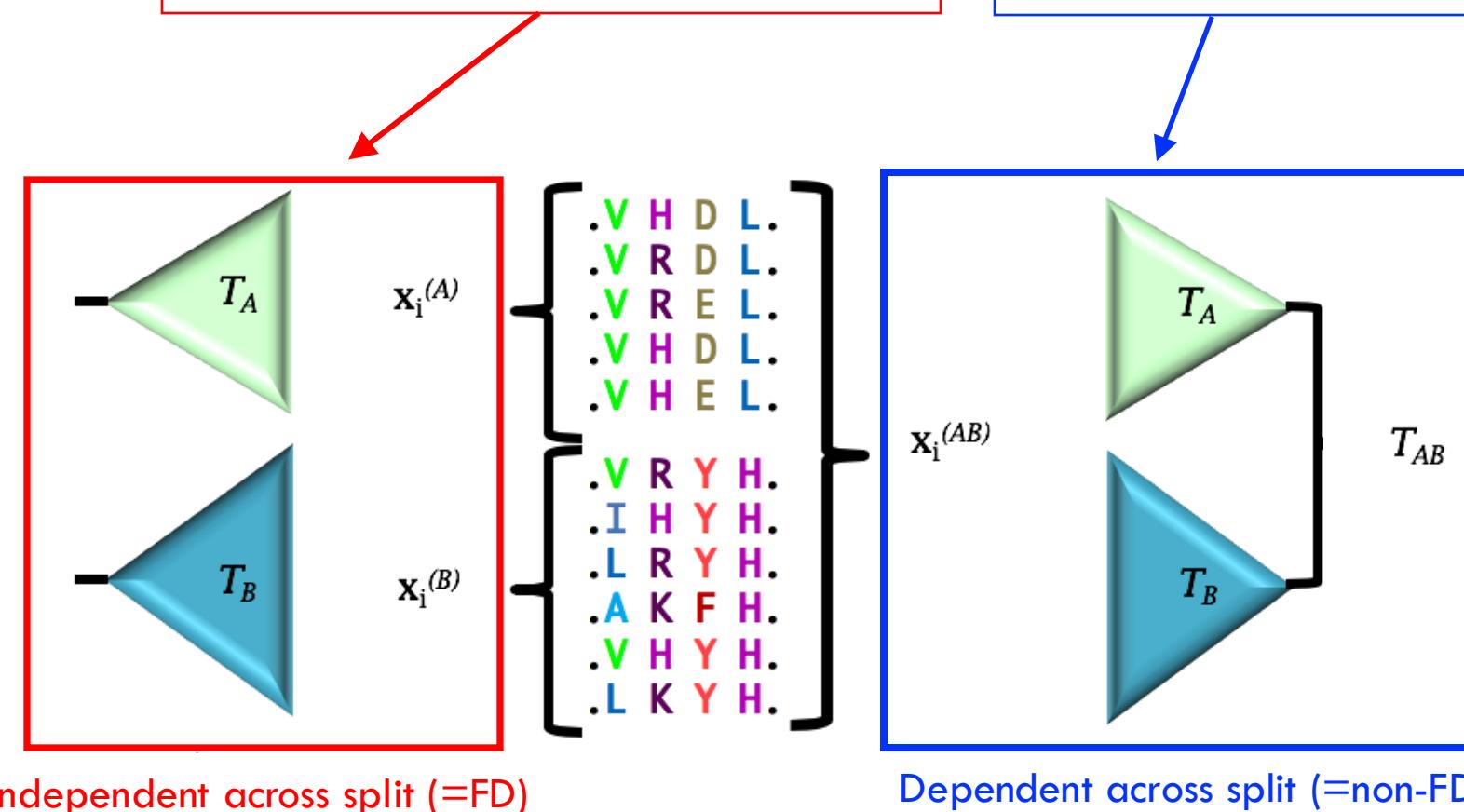


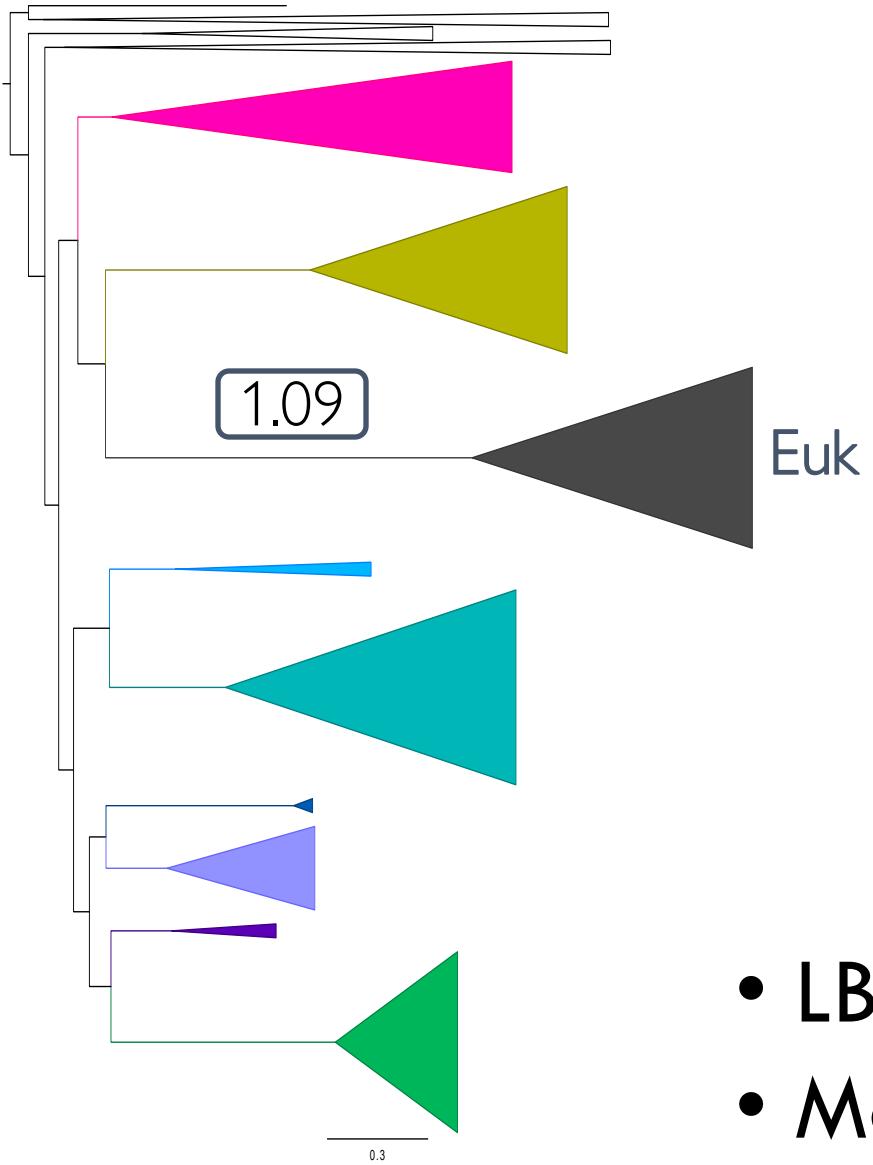
FunDi : identifies FD sites along a specific branch taking into account the phylogeny (ML framework)

FunDi mixture model

- For each site, FunDi allows for FD and non-FD evolution across a pre-specified split:

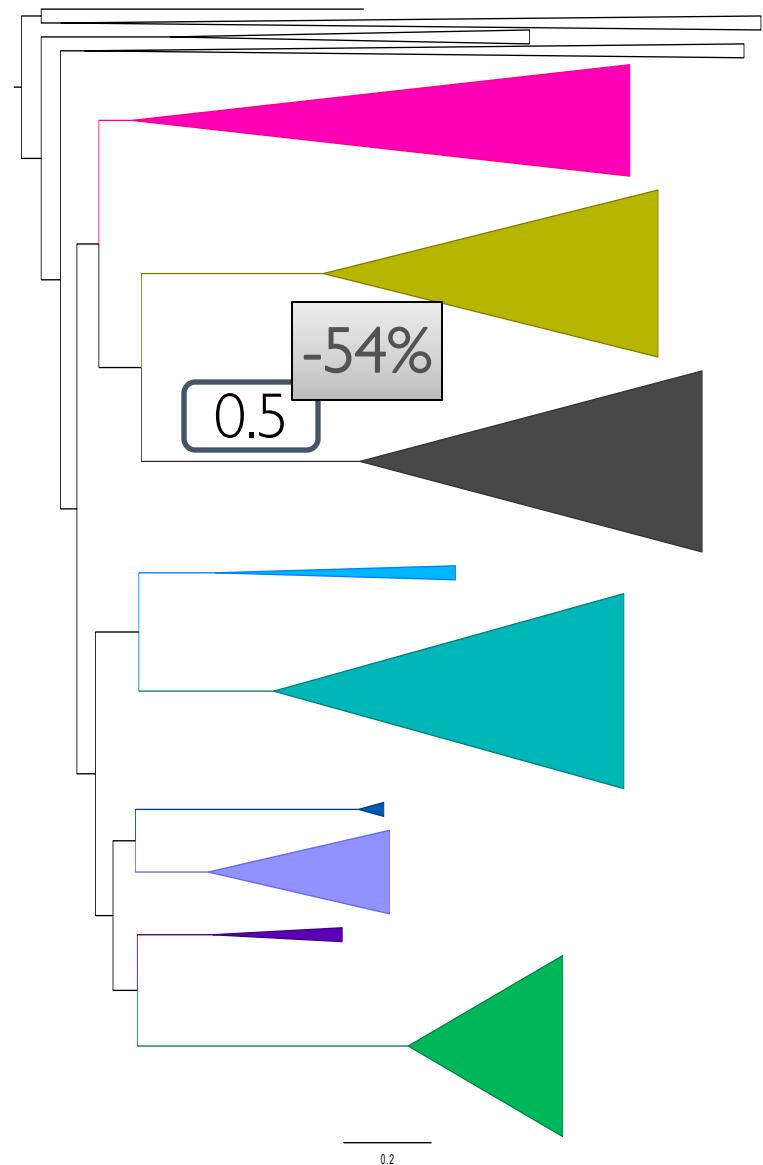
$$L_i(\theta, T_{AB}) = \boxed{\rho (P(\mathbf{x}_i^{(A)}; \theta, T_A)) (P(\mathbf{x}_i^{(B)}; \theta, T_B))} + \boxed{(1 - \rho) (P(\mathbf{x}_i; \theta, T_{AB}))}$$





- LBA
- Molecular clock

Euk
- 'FD sites'



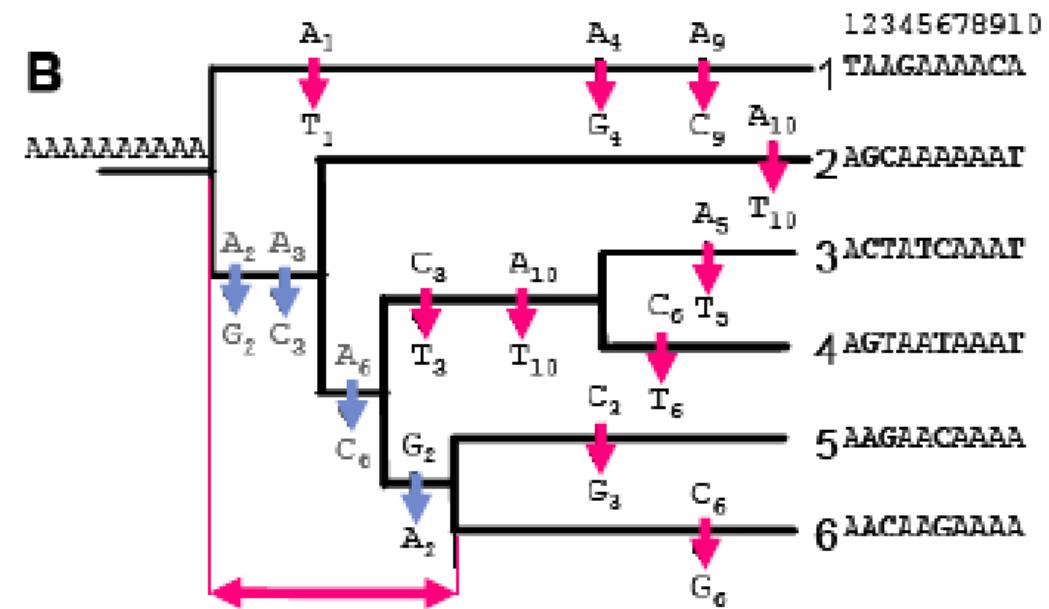
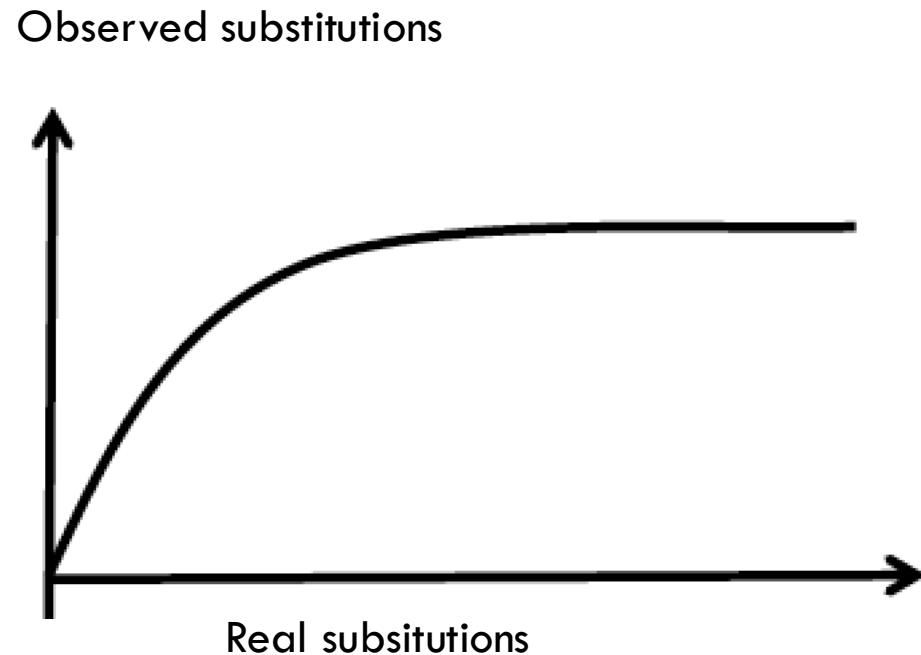
Break

PART 2

Real examples of ‘deep’ phylogenetic problems and how we tried to address them

Single gene trees are not enough to resolve 'ancient relationships'

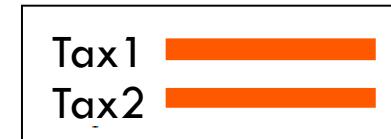
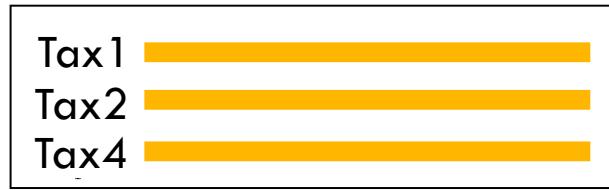
"Ancient" signal erased by more recent substitutions



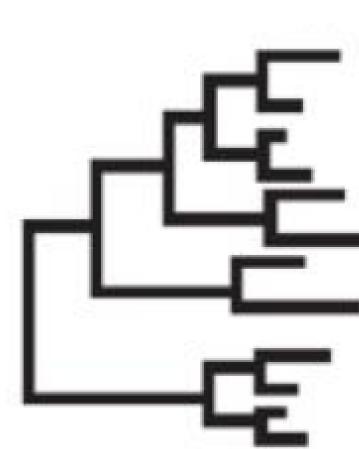
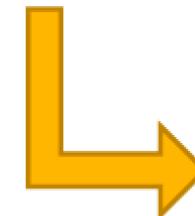
Supermatrices

Combine weak phylogenetic (historical) signal from many genes

Attenuate individual bias (IF RANDOM)



CHECK FOR CONGRUENCE



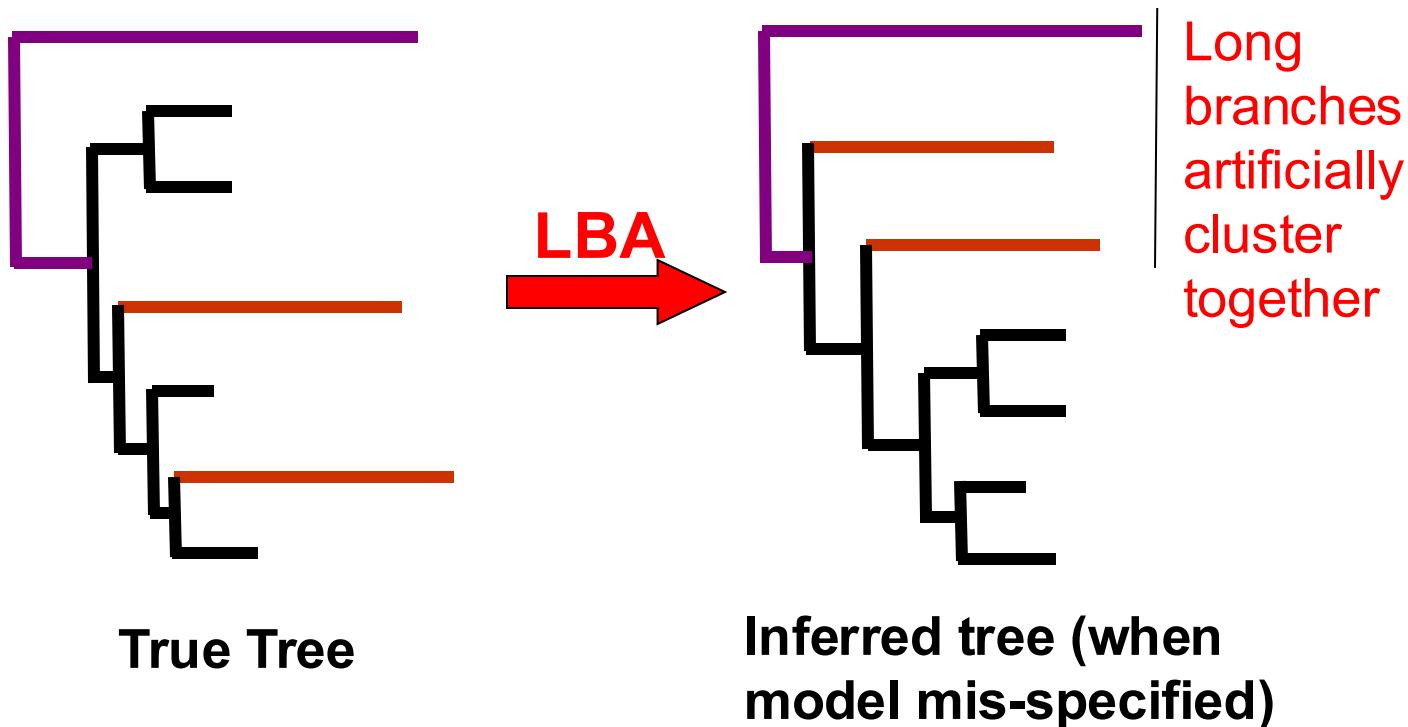
What can affect your topology

- Taxon sampling
 - Long branching taxa
 - Taxa with compositional bias
 - Contaminated data
- Gene/site sampling
 - Heterotachy
 - Saturated sites
- Model misspecification
 - LBA
- Highways of HGT
 - Consistently conflicting with vertical signal
- (many other things...)

Example 1: Effect of model misspecification and of fast-evolving sites

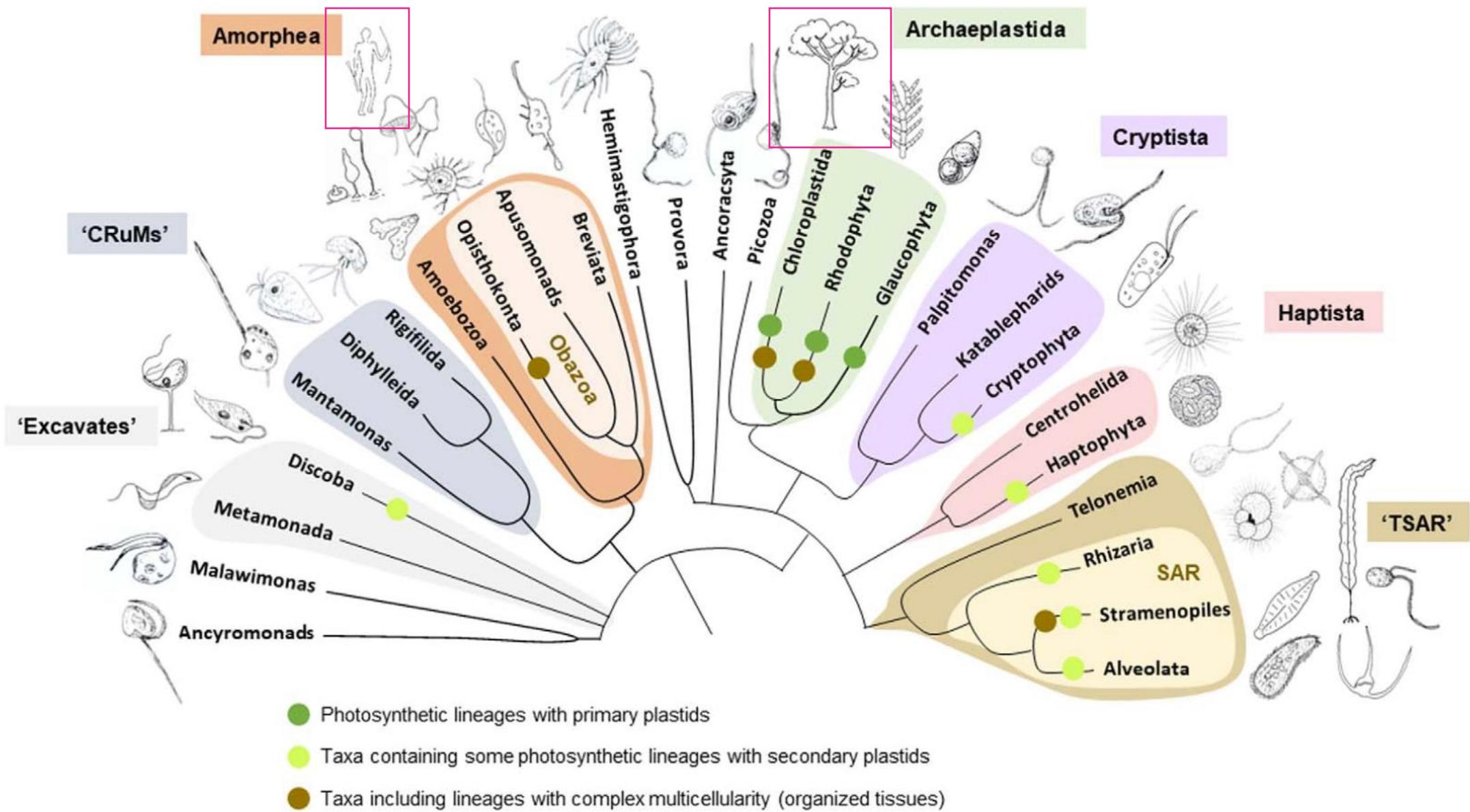
Model misspecification: statistical inconsistency

Long Branch Attraction (LBA) Artefact

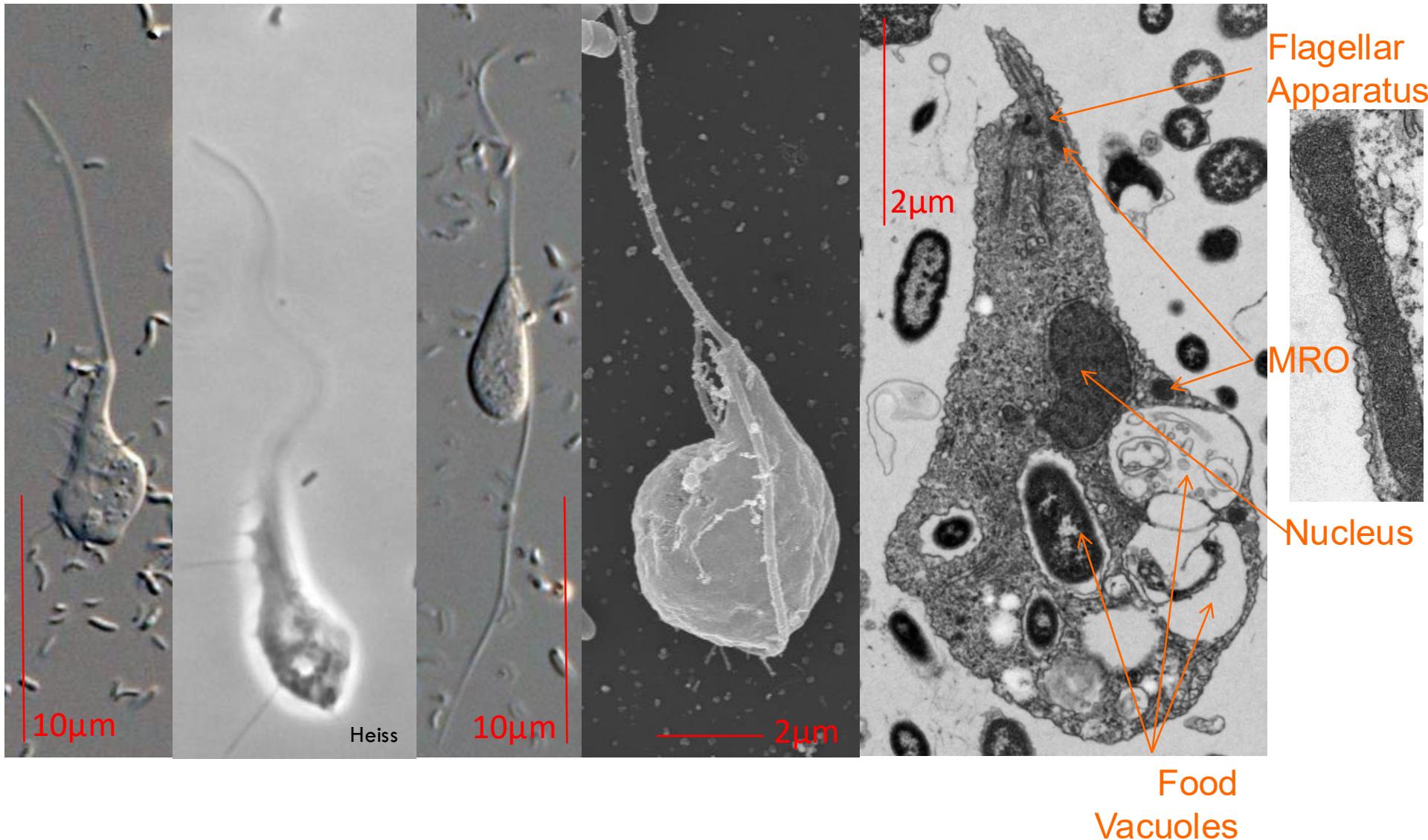


Adding more data *strengthens* artefact
→ statistical inconsistency

Tree of eukaryotes



Pygsuia biforma (Brown,...Roger 2013)



Two different topologies within Obazoa are supported by different phylogenetic models



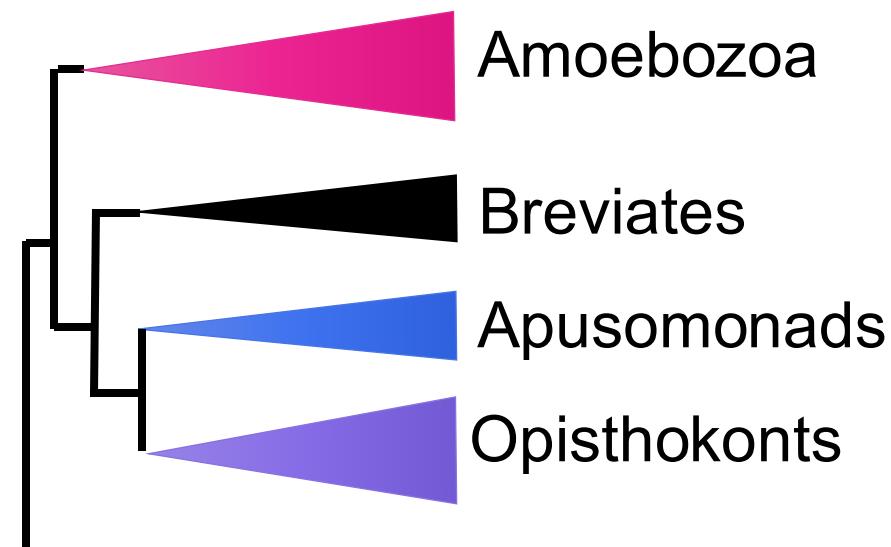
Opisto + Breviates + Apusomonads = OB_Azoa

Two different topologies within Obazoa are supported by different phylogenetic models

ML-BS = 98%



Bayes posterior prob. = 1.0



ML – LG+ Γ
Bayes – LG+ Γ

Bayes – CAT-Poisson+ Γ
Bayes – CAT-GTR+ Γ

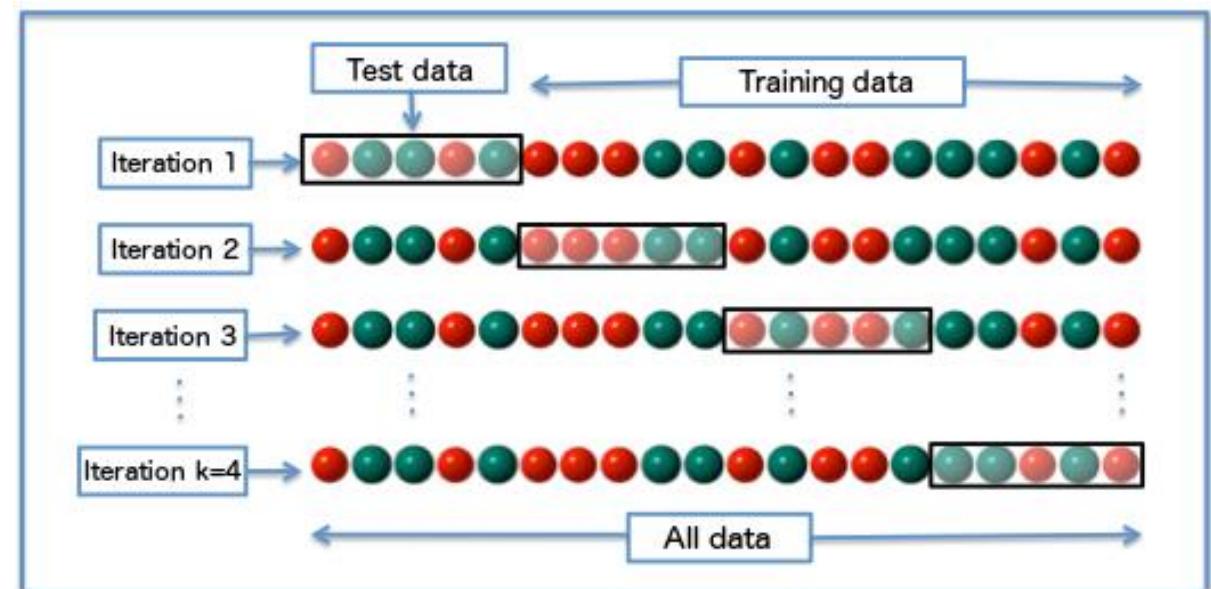
How to decide which is real and which is artefact?

- One of two topologies is likely artefactual resulting from mis-specified model
- Test which substitution model fits better
 - E.g., Cross-validation, Bayes factors, Posterior prediction

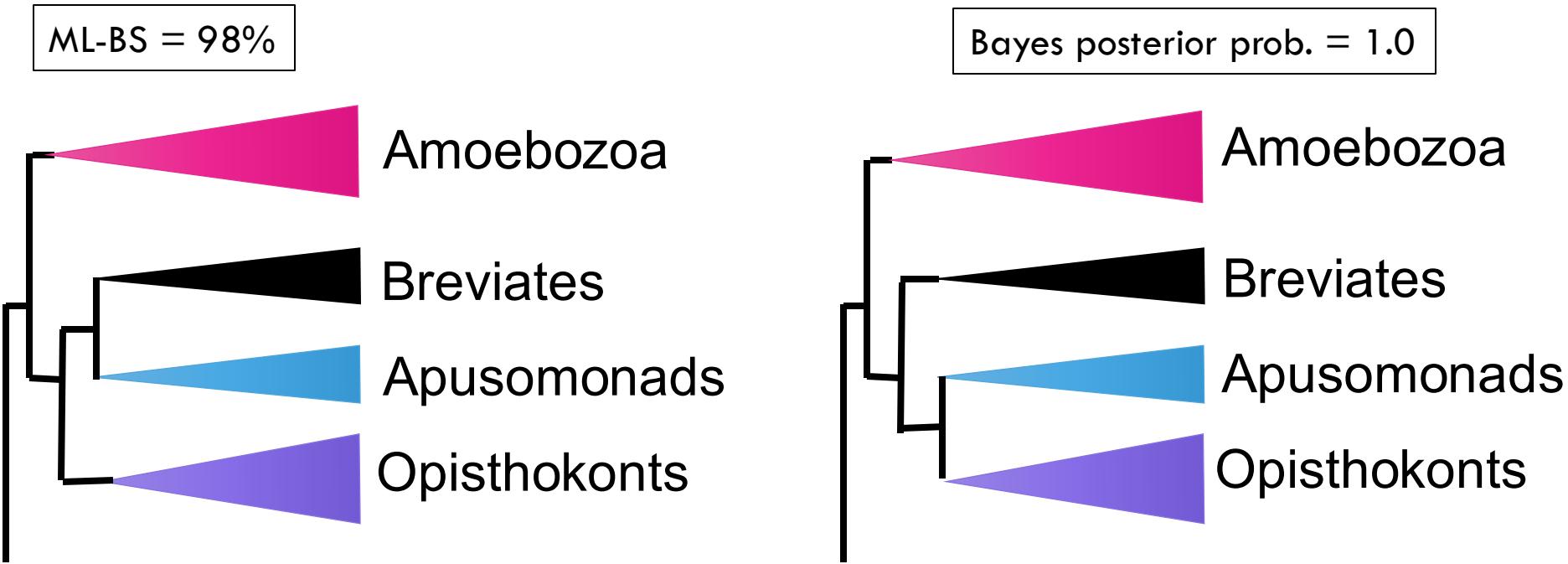
How to decide which is real and which is artefact?

Cross-validation:

- 1) parameters of the model estimated on the learning set
- 2) these parameter values are then used to compute the likelihood of the test set = **how well the test set is 'predicted' by the model?**
- 3) Repeat over all partitions and average the likelihood
- 4) Repeat for each model and compare



Cross-validation favors CAT-GTR over LG



ML – LG+ Γ
Bayes – LG+ Γ

Bayes – CAT-Poisson+ Γ
Bayes – CAT-GTR+ Γ



Cross validation

How do decide which is real and which is artefact?

- One of two topologies is likely artefactual resulting from misspecified model
- Test which substitution model fits better
 - Cross-validation
- Try to eliminate ‘noisiest’ data
 - Fast-evolving site removal
 - Fast-evolving gene removal
 - Fast-evolving taxon removal
 - Recoding

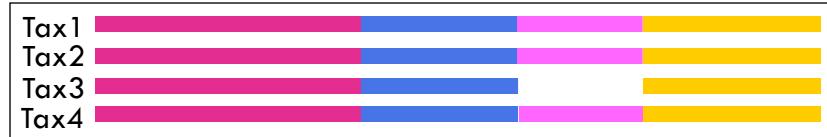
Removal of fast-evolving sites

Goal: can we yield the same topology
under LG+G as under CAT+GTR
after we remove poorly modelled sites?

Fast Evolving Sites removal

Fast-evolving sites : carry the ‘noisiest’ signal (most saturated sites)

Initial alignment



Tax1

Tax2

Tax3

Tax4

Iteration 1

Tax1

Tax2

Tax3

Tax4



Reconstruct tree

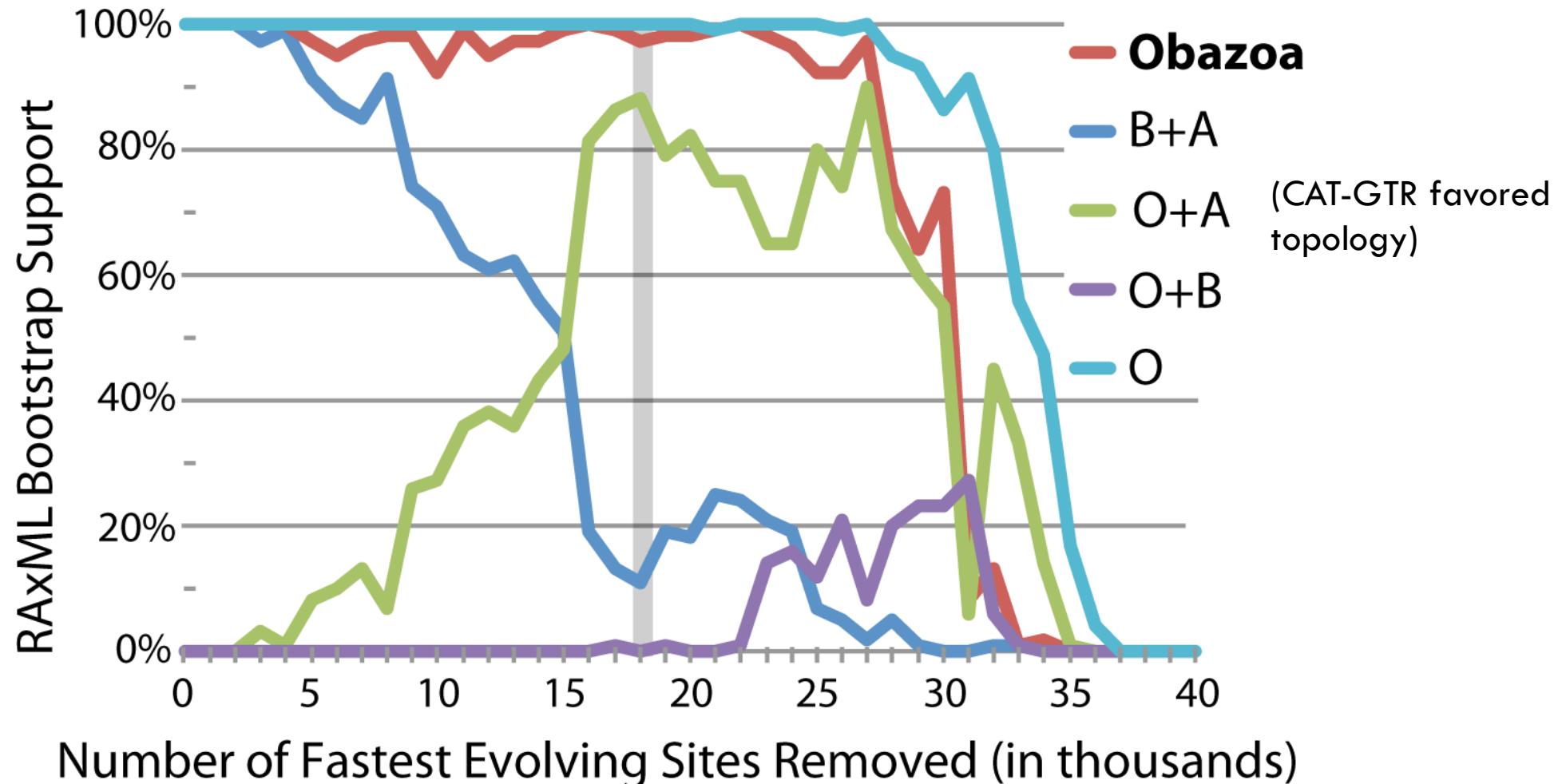
→ 40 steps of removal of 1000 sites (evolutionary rates estimated by IQTREE for example)

Step	# sites left	
1	43615	Tree 1
2	42615	Tree 2
...		
4	40615	Tree 4
...		
40	3615	Tree 40

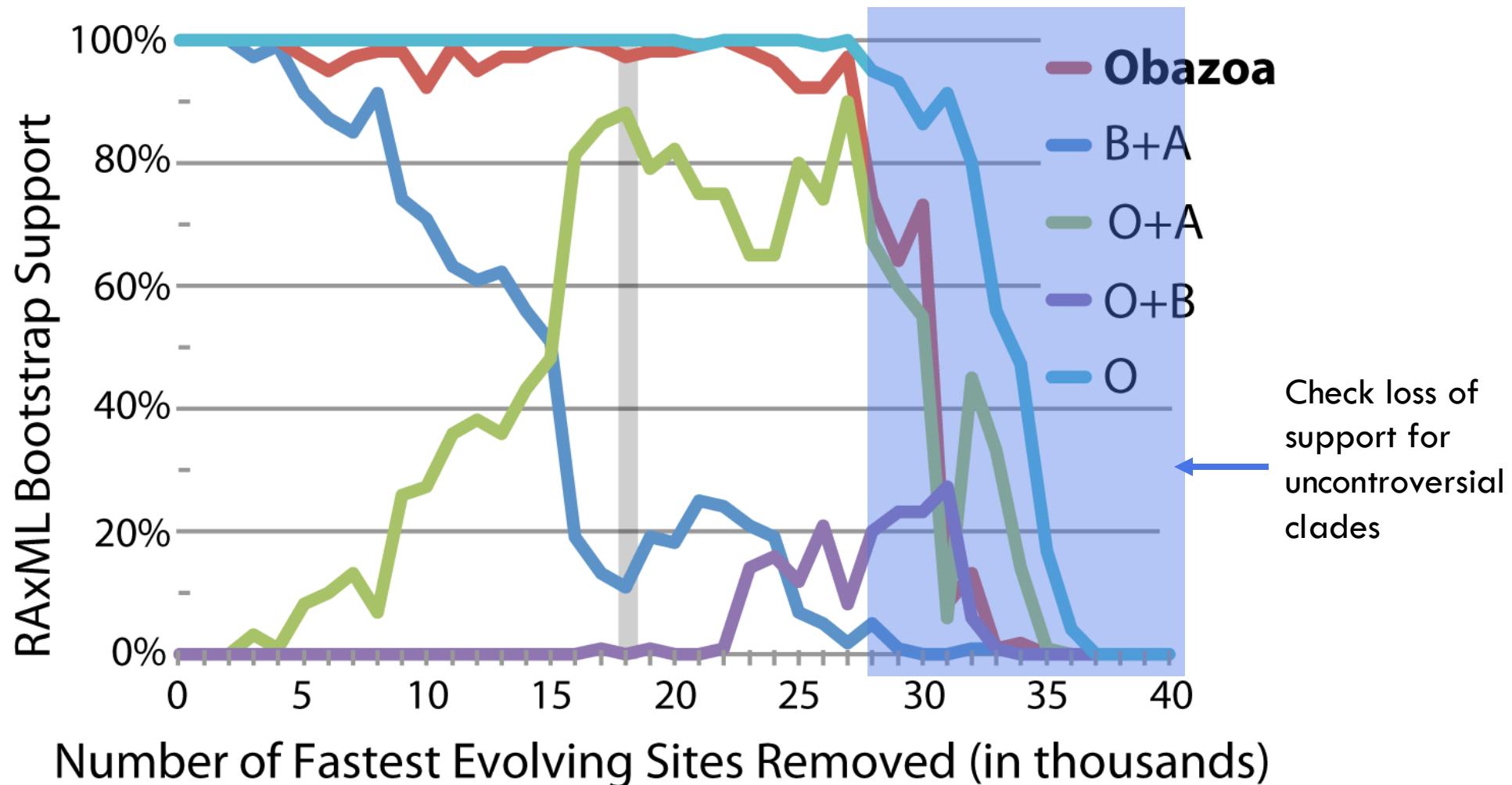
Estimate support for conflictual clades as we remove Fast-Evolving Sites (under LG+G)



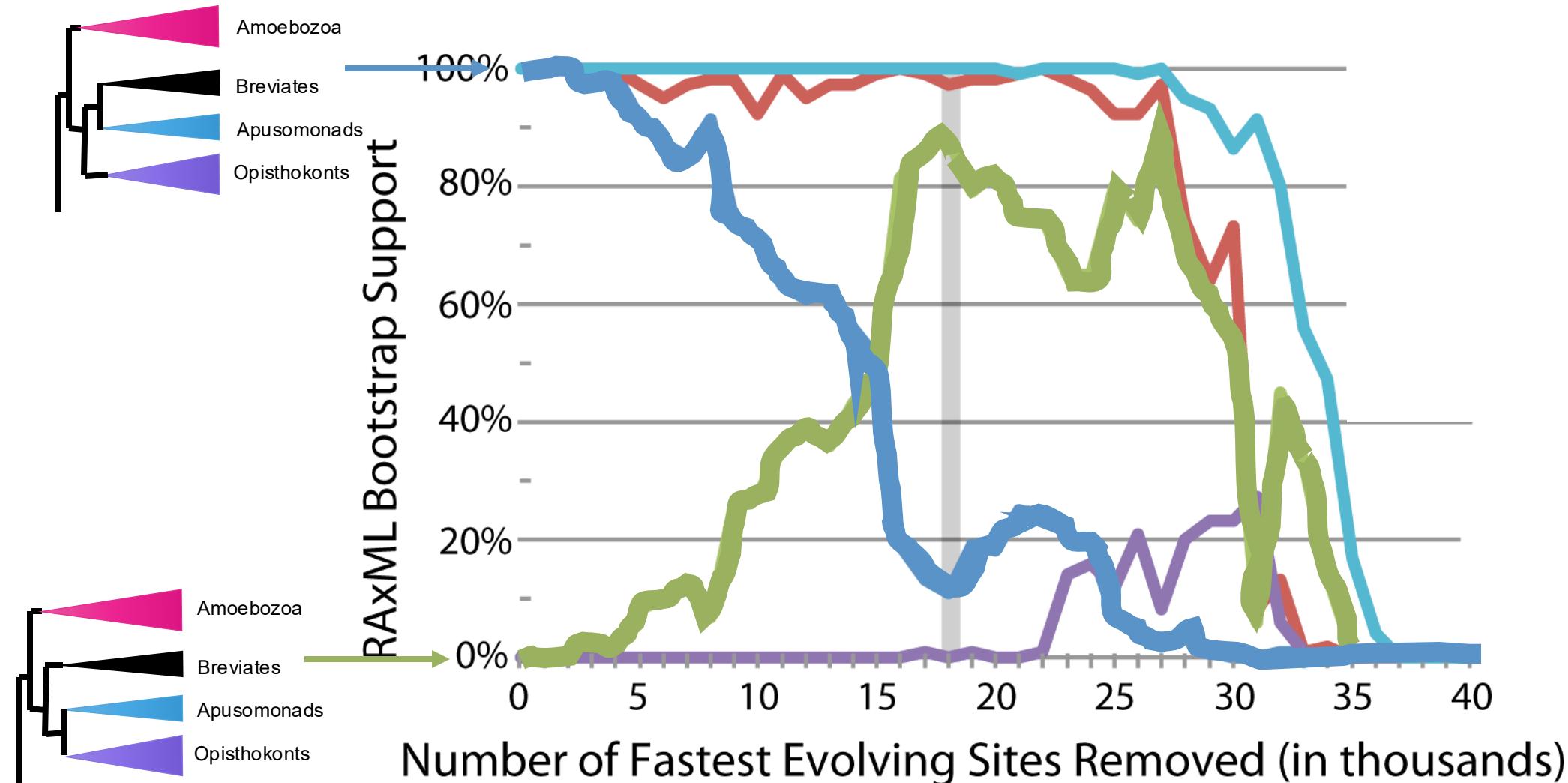
Estimate support for conflictual clades as we remove Fast-Evolving Sites (under LG+G)

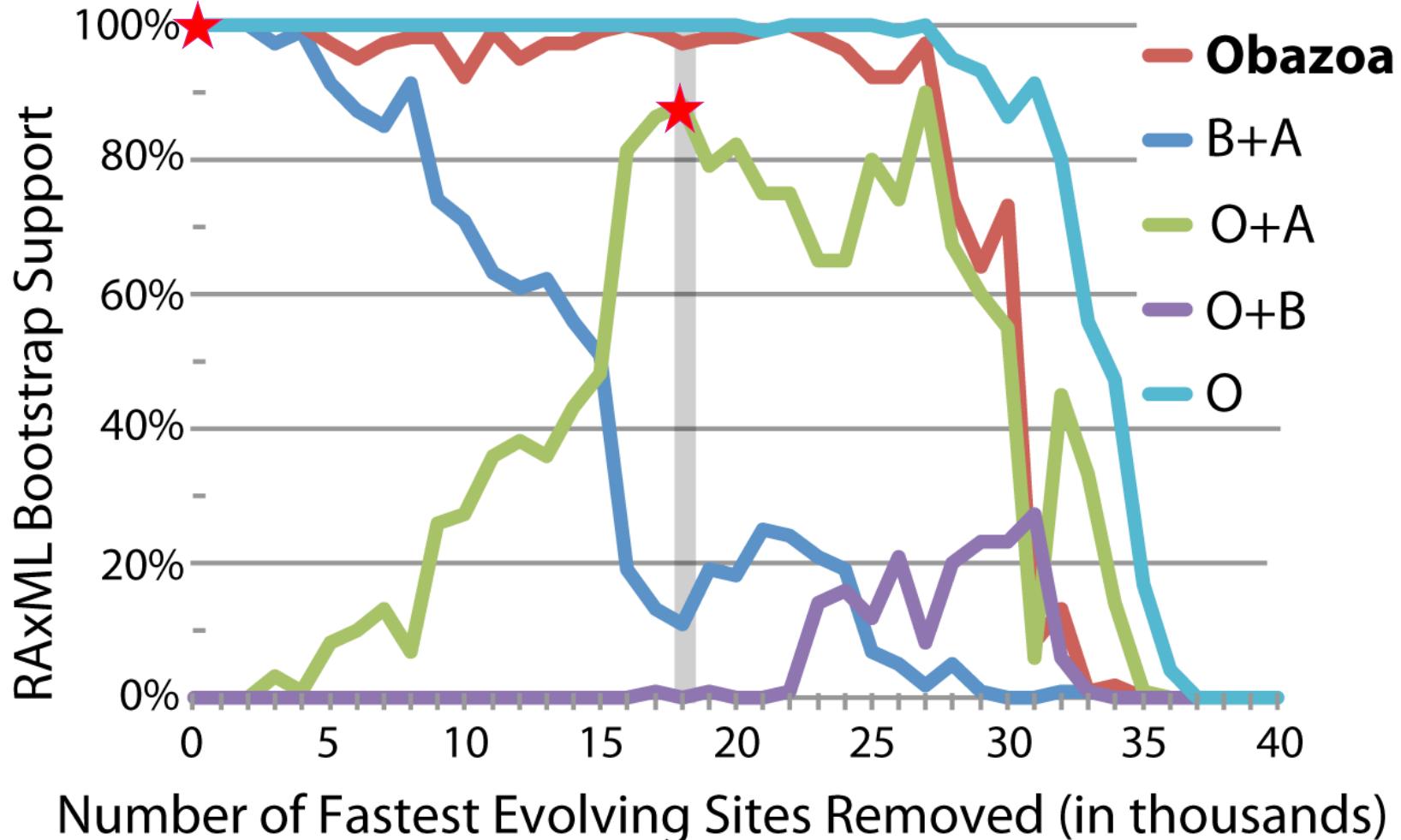


Estimate support for conflictual clades as we remove Fast-Evolving Sites (under LG+G)



Breviates+Apusomonads (B+A) topology vs. Apusomonads+Opisthokonts (O+A)







Removal of 18,000 fastest-evolving sites



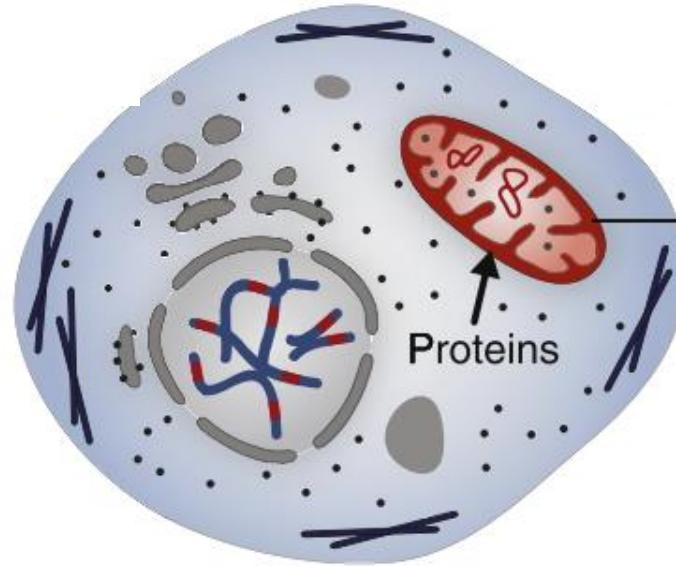
Fast-evolving site removal: warning

- Poor proxy for heterotachous site removal
 - Fast sites in themselves are not necessarily a problem if they are fast across the entire tree
- In practice, fast sites seem to overlap to some extent with sites whose rate varies across the tree and are improperly modelled by most widely used models.
- You also remove the most saturated sites, which are usually poorly modelled

**Example 2: Effect of amino-acid
preference change over the tree**

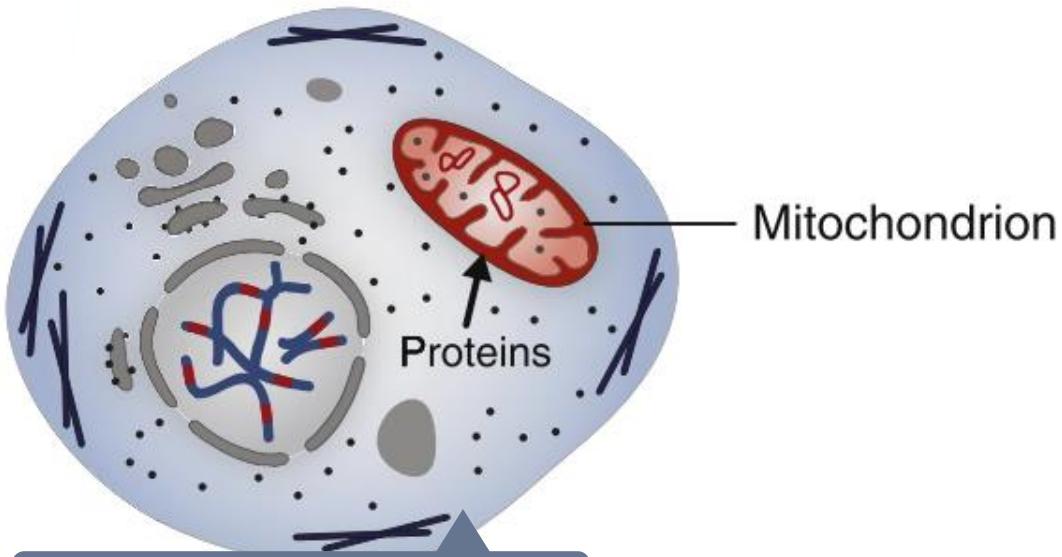
A brief account of the little we know about the origin of eukaryotes





Complex eukaryotic cell

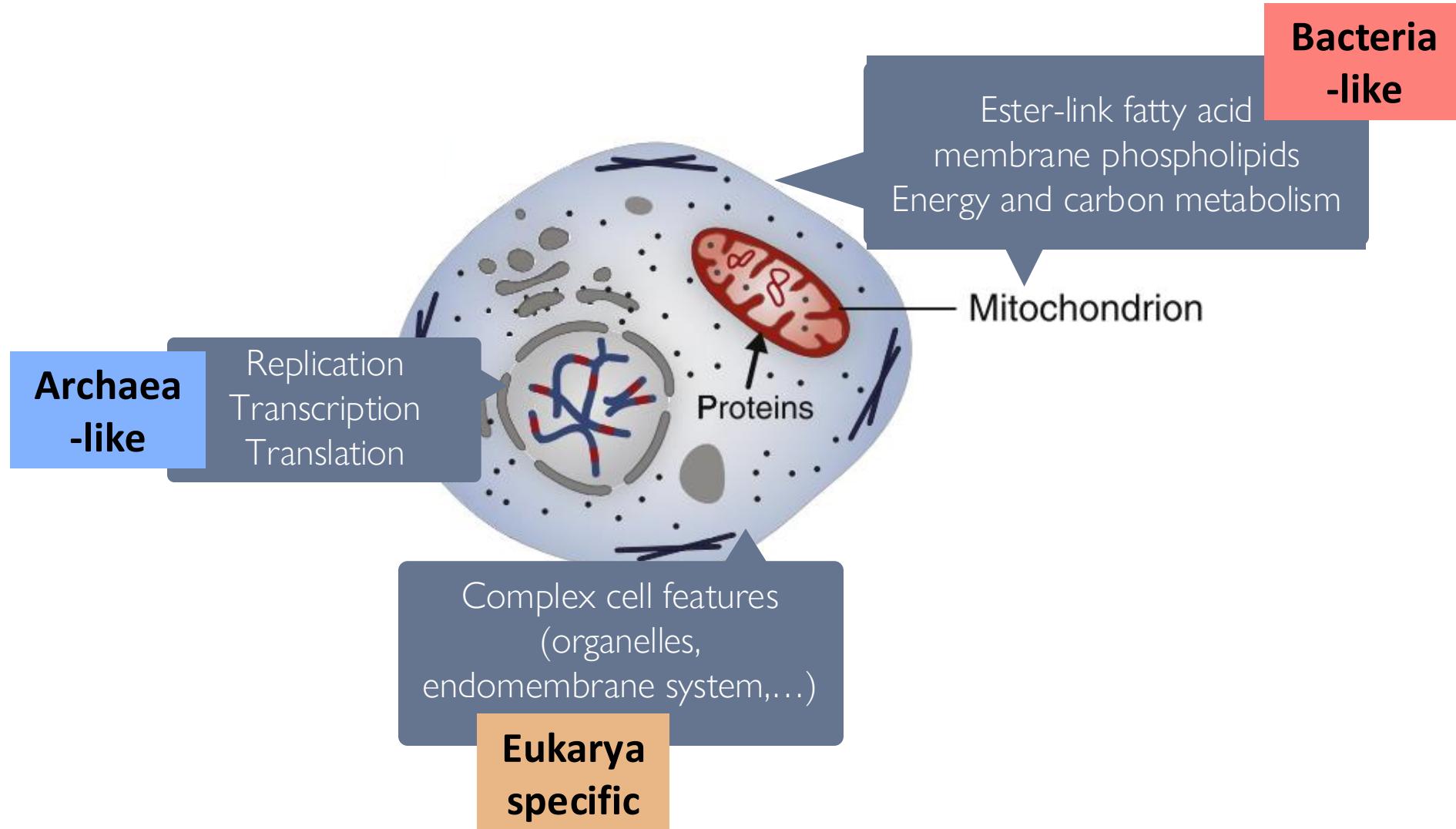
The chimeric nature of eukaryotes



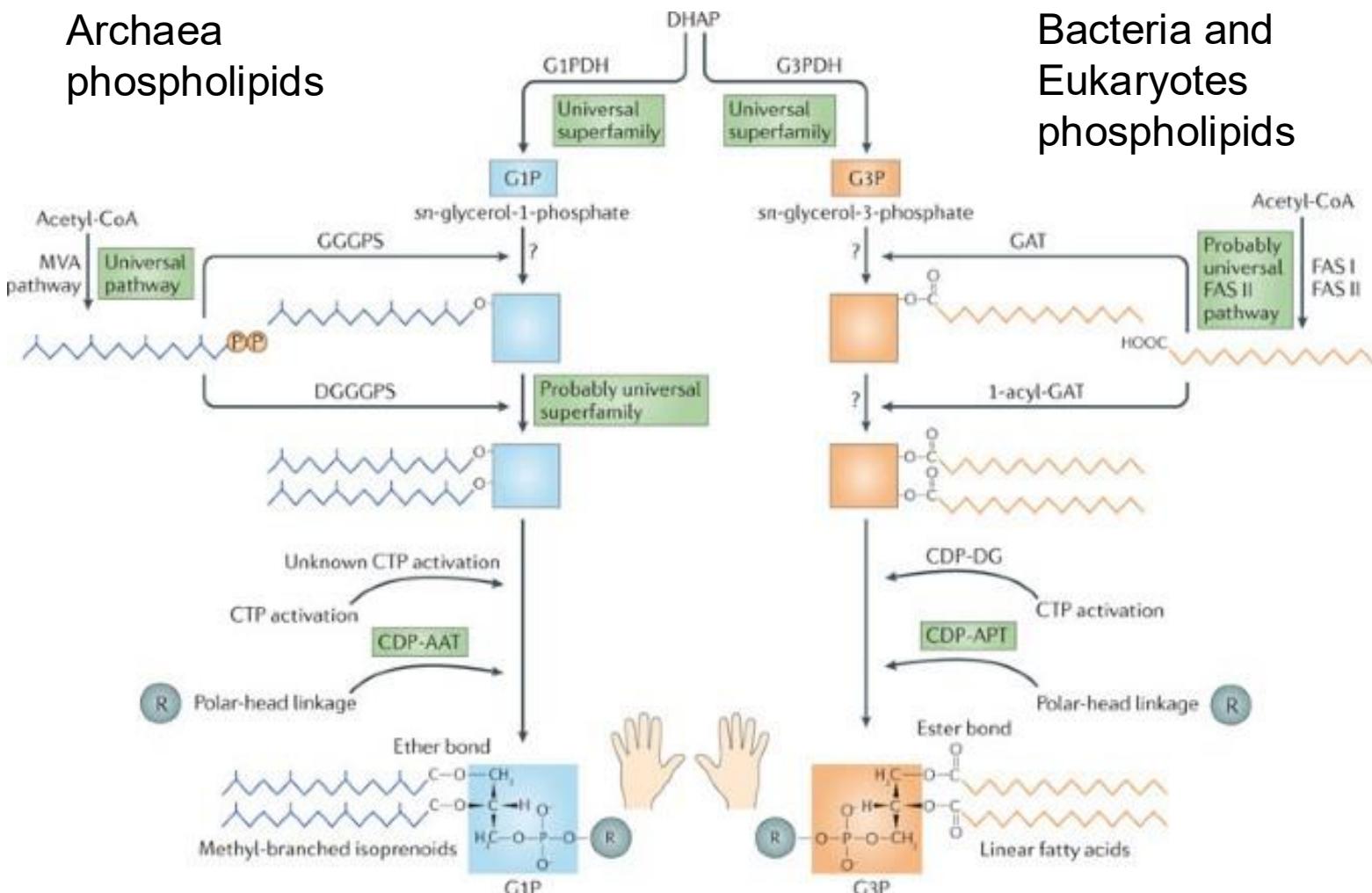
Complex cell features
(organelles,
endomembrane system,...)

**Eukarya
specific**

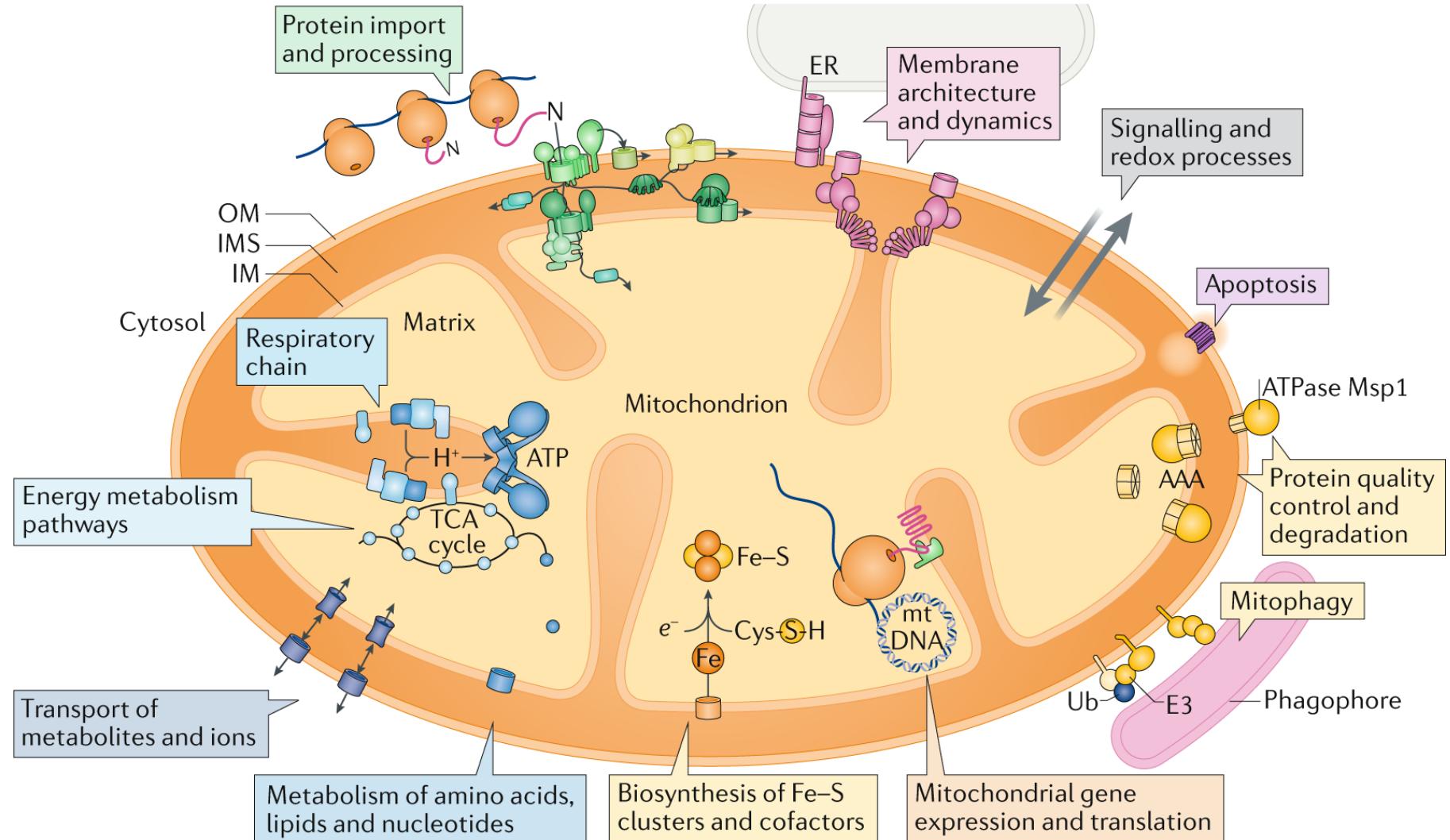
The chimeric nature of eukaryotes



Eukaryotic lipids resemble bacterial ones



Mitochondria have diverse and crucial metabolic roles for the eukaryotic cell

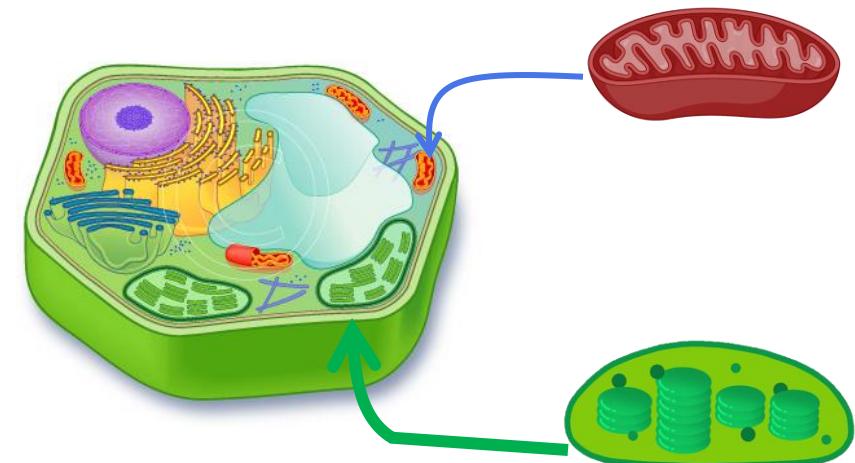
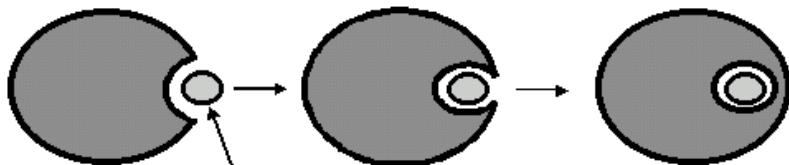


The notion that eukaryotes are chimeric in nature is not new



Lynn Margulis
(1938-2011)

- Endosymbiont theory: Mitochondria and chloroplasts were once free-living bacteria
- First proposed by Altmann (1890) and Mereschkowsky (1905), developed into modern form by Lynn Margulis

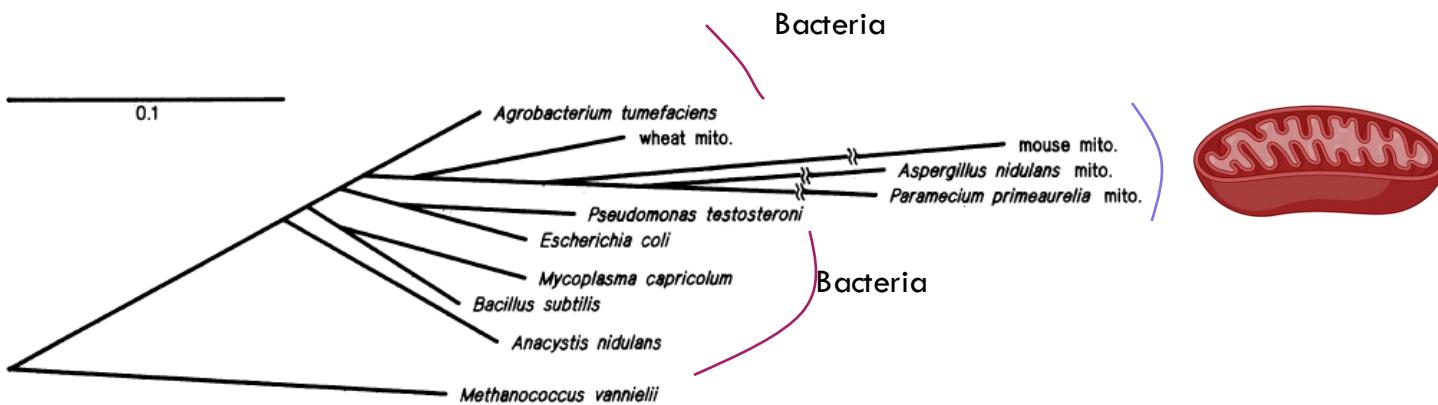


The notion that eukaryotes are chimeric in nature is not new

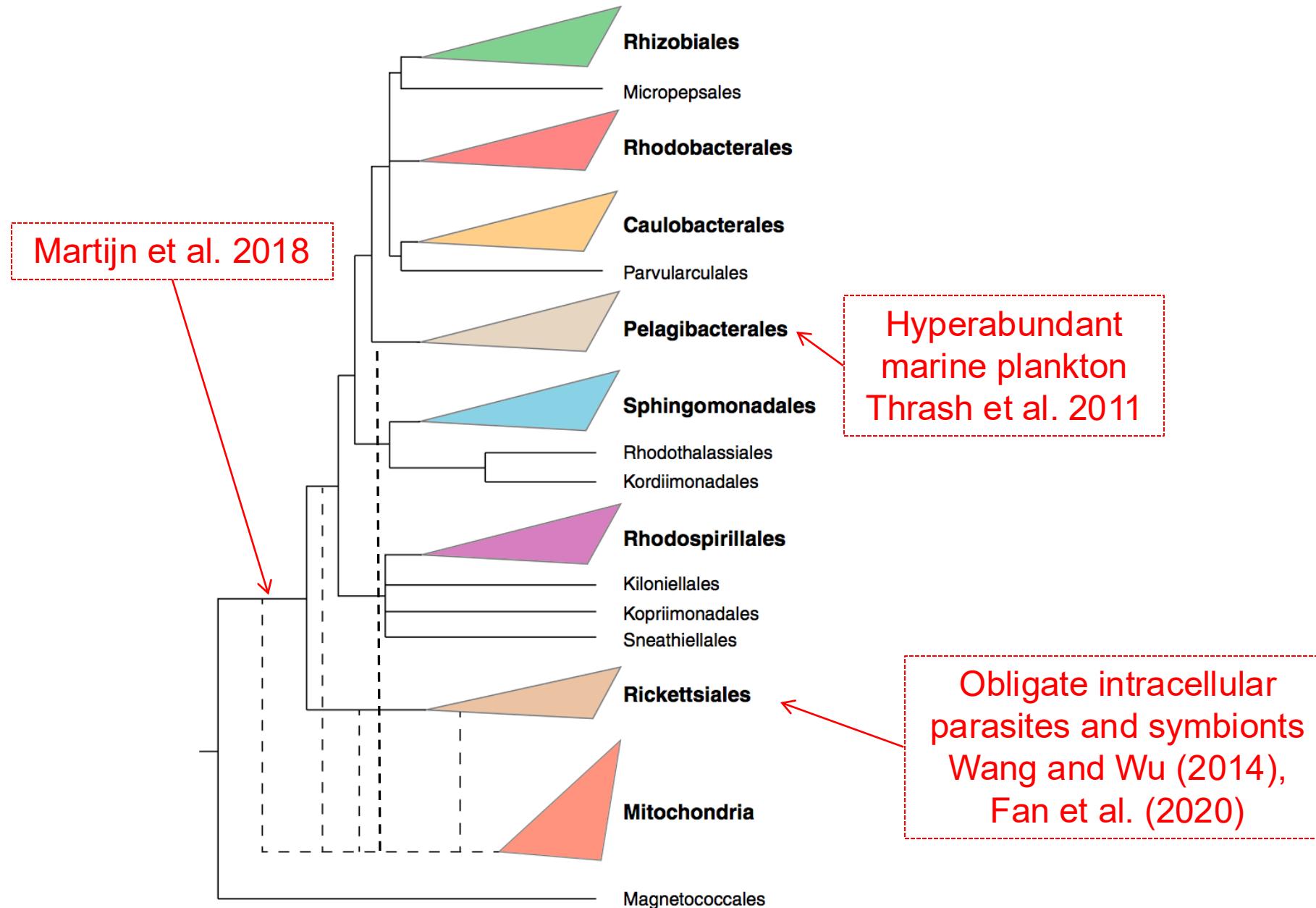


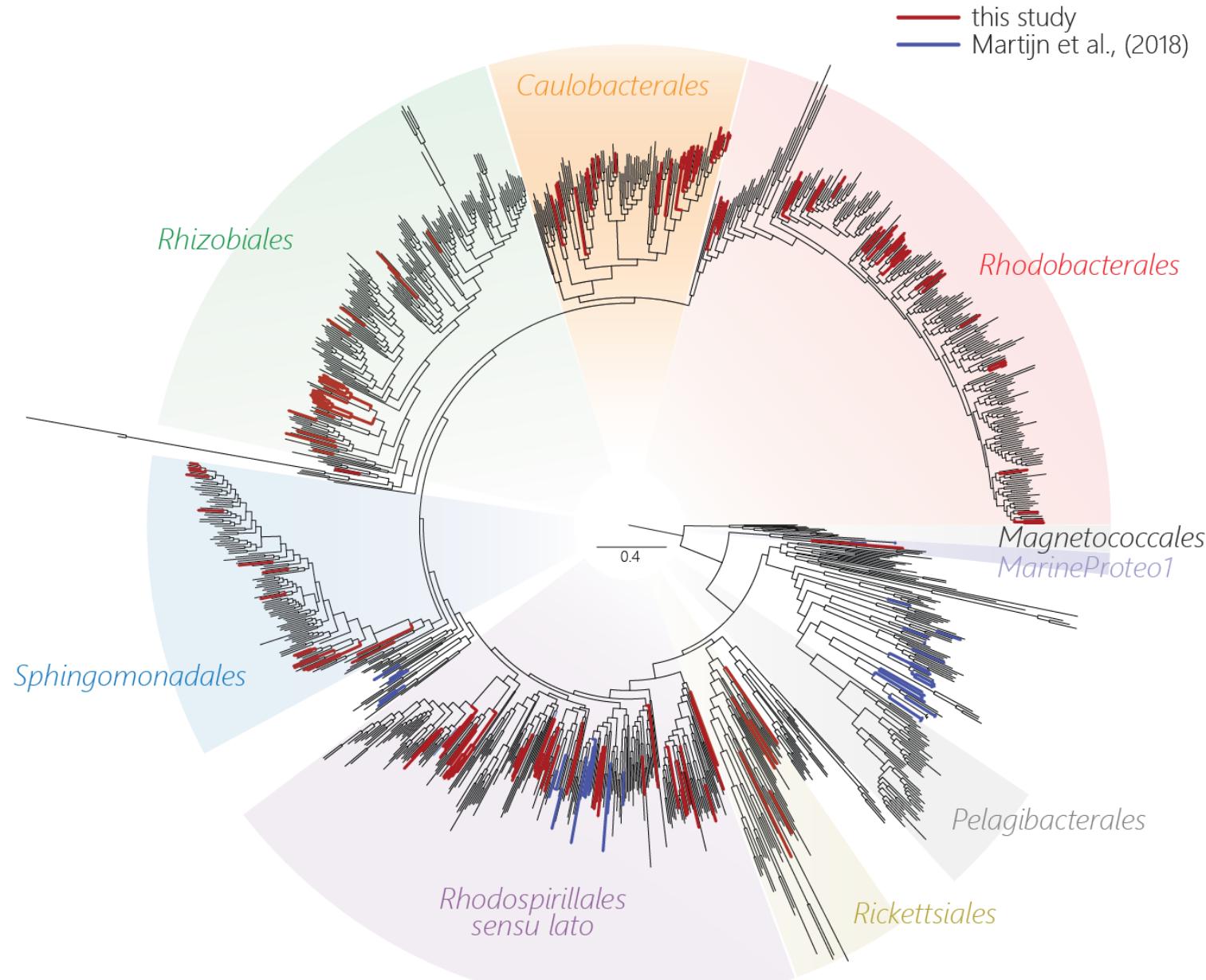
Lynn Margulis
(1938-2011)

- Endosymbiont theory: Mitochondria and chloroplasts were once free-living bacteria
- First proposed by Altmann (1890) and Mereschkowsky (1905), developed into modern form by Lynn Margulis



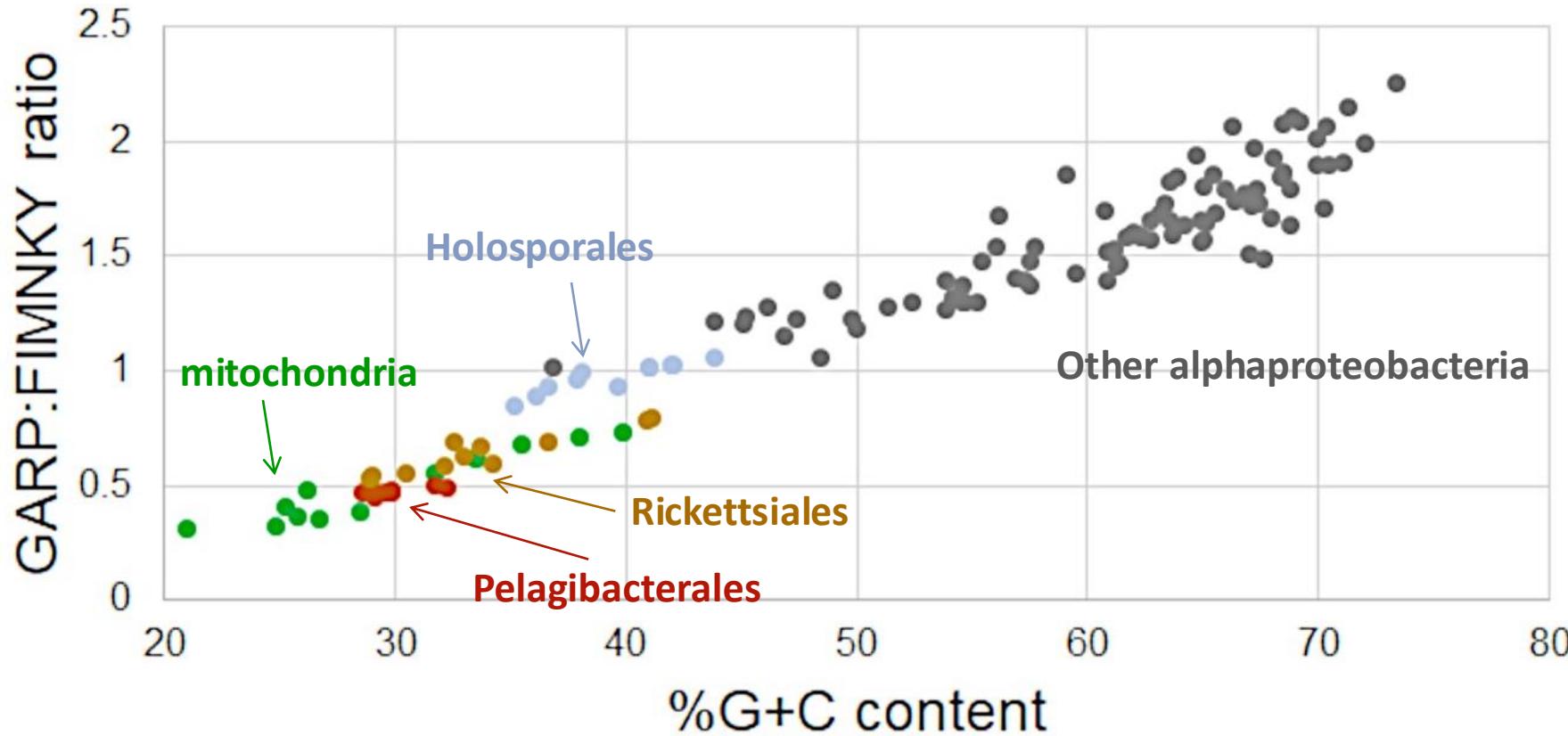
Controversy over position of mitochondria within Alphaproteobacteria





- More markers:
108 proteins of mitochondrial origin
- More taxa:
150 non-marine alphaproteobacterial MAGs
(microbial mats, microbialites and lake sediments)
- New model:
Gmix phylogenetic model

The proteome AA composition varies with the genomic GC content



A site-and-branch-heterogeneous profile mixture model GFmix



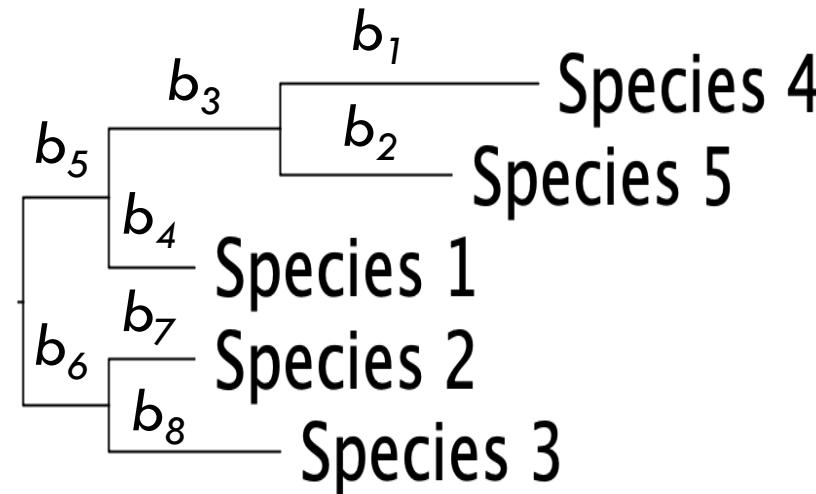
Edward Susko

For each branch assume there's a specific ratio of frequencies of GARP:FYMINK

- call this the 'b' parameter

$$b_x = \frac{f_{G,A,R,P}}{f_{F,I,M,N,K,Y}}$$

Then for the tree we have:

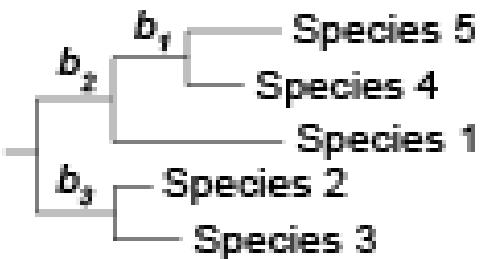


GFmix modifies the site profile mixture classes for each branch based on the b parameter

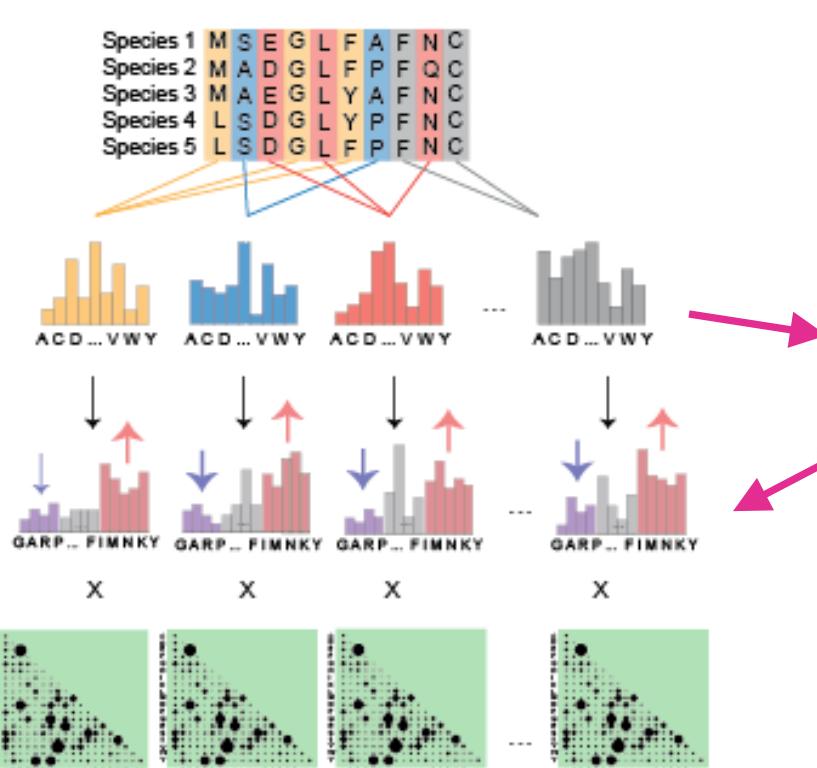
A site-and-branch-heterogeneous profile mixture model GFmix



Edward Susko



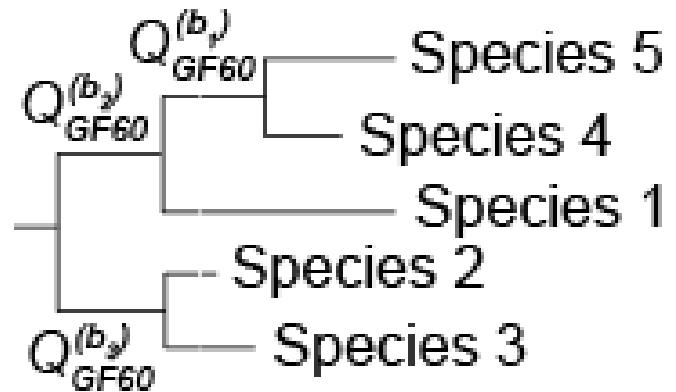
Species 1	M	S	E	G	L	F	A	F	N	C
Species 2	M	A	D	G	L	F	P	F	Q	C
Species 3	M	A	E	G	L	Y	A	F	N	C
Species 4	L	S	D	G	L	Y	P	F	N	C
Species 5	L	S	D	G	L	F	P	F	N	C



$$Q_{GF60}^{(b_1)} \downarrow \begin{bmatrix} Q_{gf1} \\ Q_{gf2} \\ Q_{gf3} \\ \vdots \\ Q_{gf4} \end{bmatrix}$$

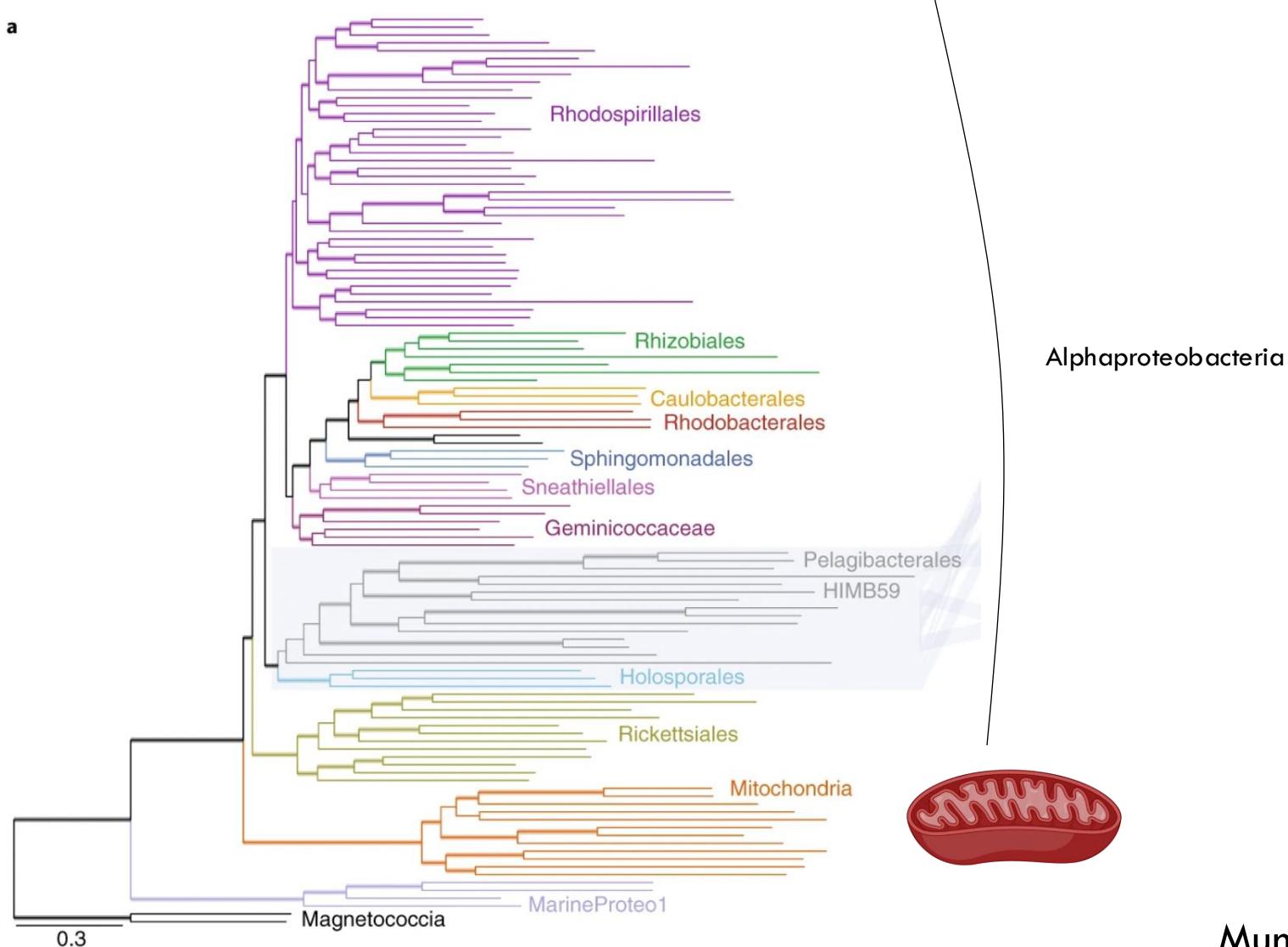
$$\vdots$$

$$Q_{GF60}^{(b_x)} \downarrow \begin{bmatrix} Q_{gf1} & Q_{gf2} & Q_{gf3} & \cdots & Q_{gf4} \end{bmatrix}$$

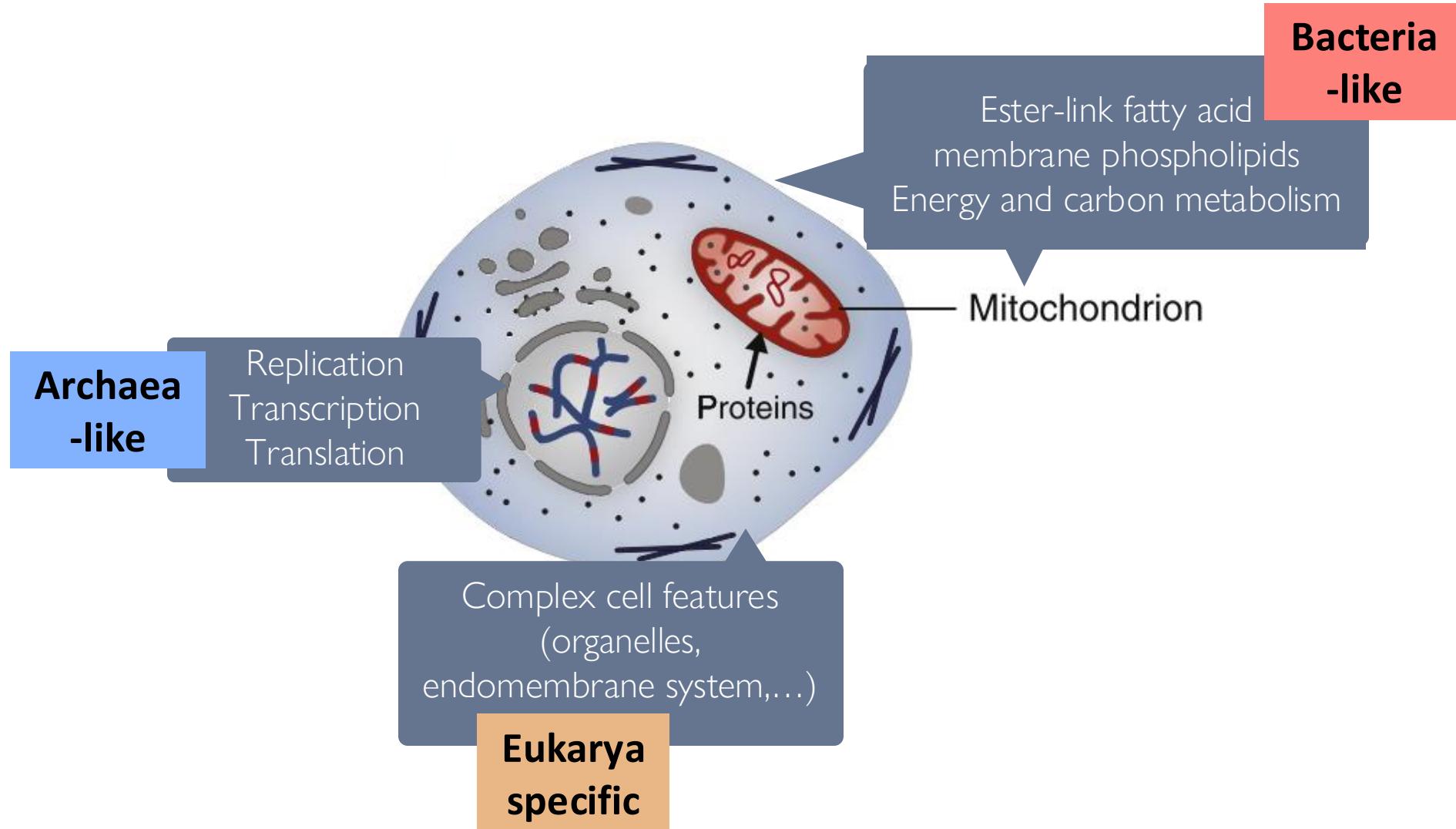


The GFmix model supports the ancestry of mitochondria from outside known diversity of alphaproteobacteria

a

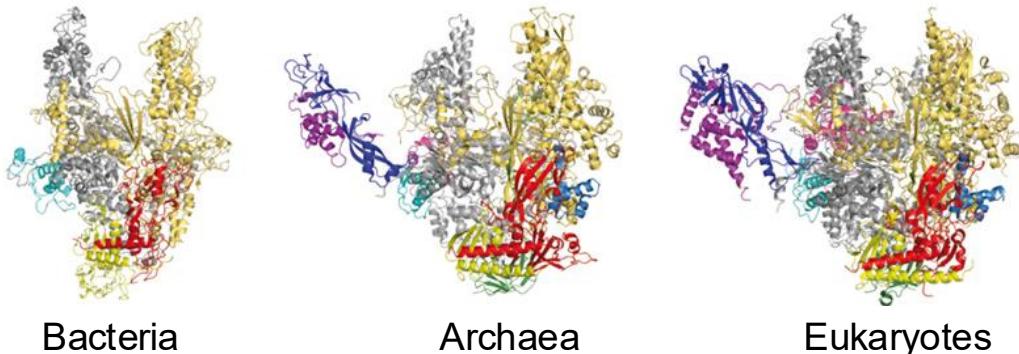


The chimeric nature of eukaryotes



Informational systems suggest an archaeal connection

Transmission and expression of genetic information show a higher similarity between eukaryotes and Archaea than with Bacteria

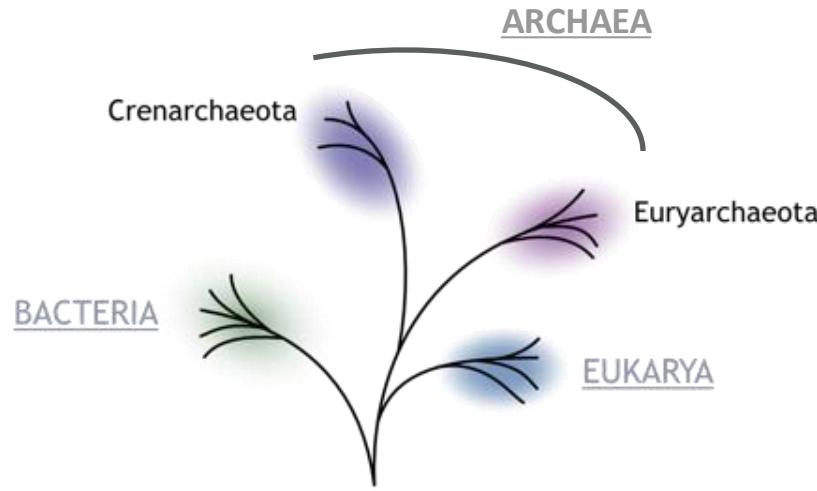


Overall architecture of RNA polymerases (RNAPs)

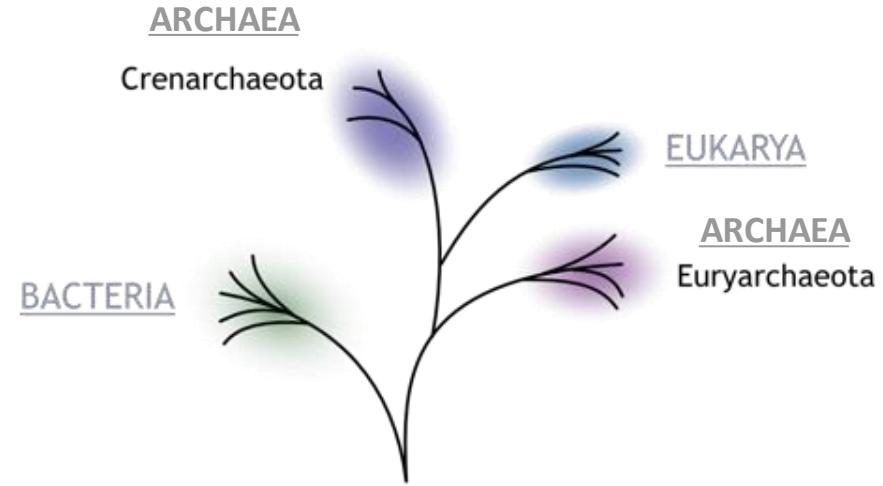
	Bacteria	Archaea	Eukaryotes
Conserved Core	β' β α αII ω	A' and A'' B' and B'' D L K	Rpb1 Rpb2 Rpb3 Rpb11 Rpb6
Archaea + Eukaryotes		H G * N P F E'	Rpb5 Rpb8 Rpb10 Rpb12 Rpb4 Rpb7
General transcription factors (GTFs)		TBP TFB TFE α TFE β /C34	Rpb9 TBP TFIIB TFIIE α TFIIE β

Subunit composition of the RNAPs

Archaea as sister-group or as ancestors of eukaryotes?



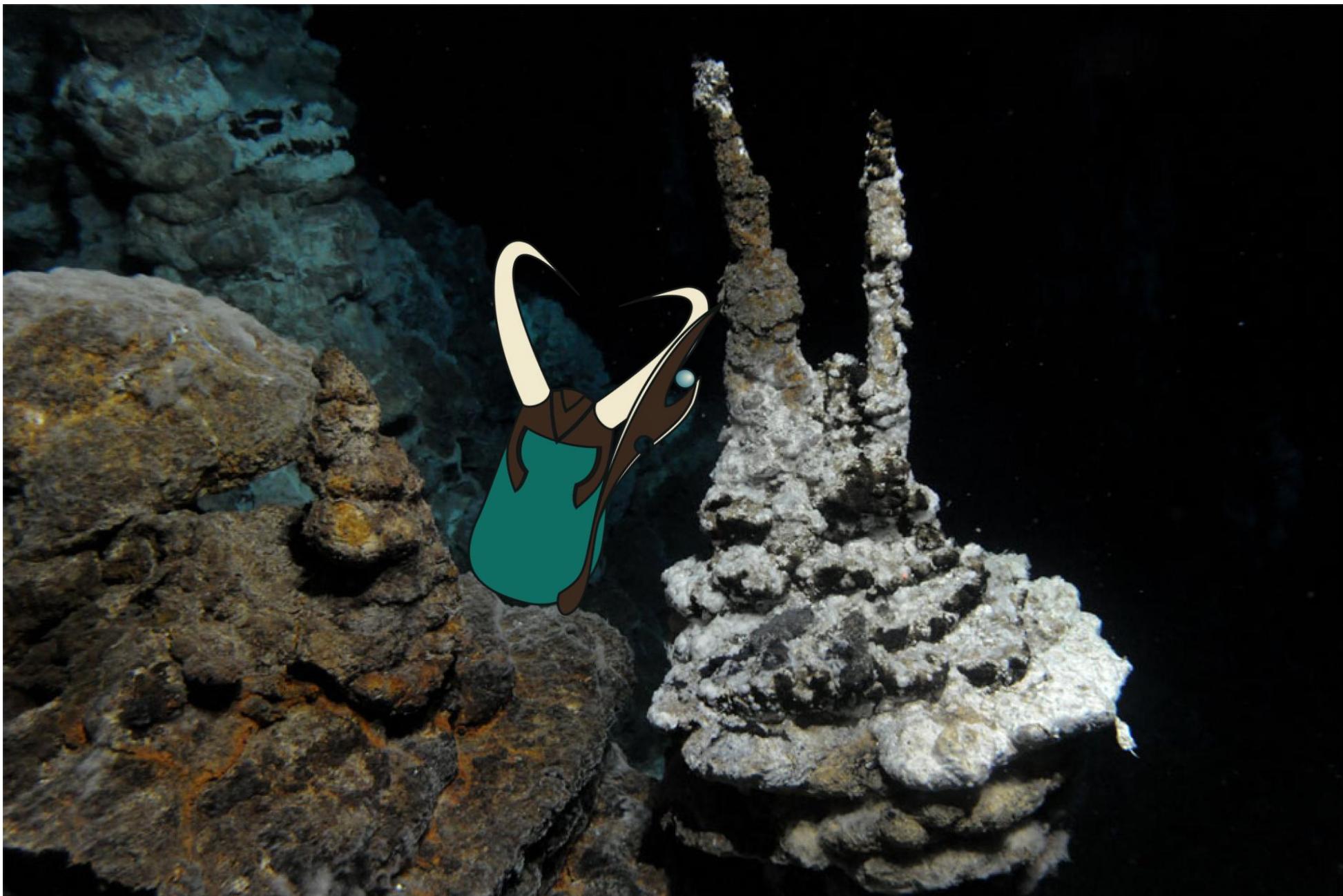
Three domain
tree of Life



Two domain
tree of Life

1990s-2000s: Phylogenetic analyses: few (informational) genes; few cultivated organisms

Culture-independent genomics (e.g. metagenomics)



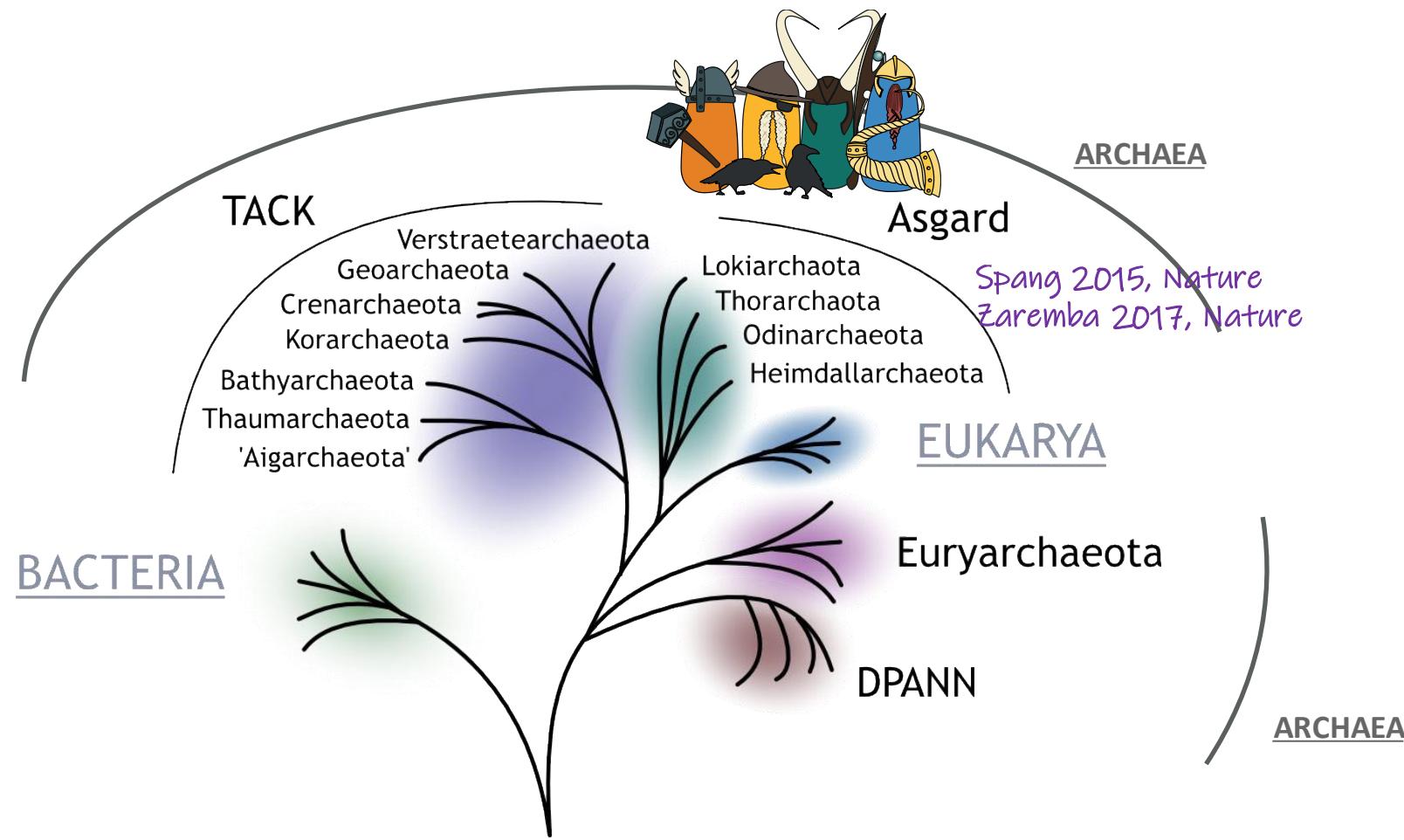
The unveiling of Asgard archaea through metagenomics



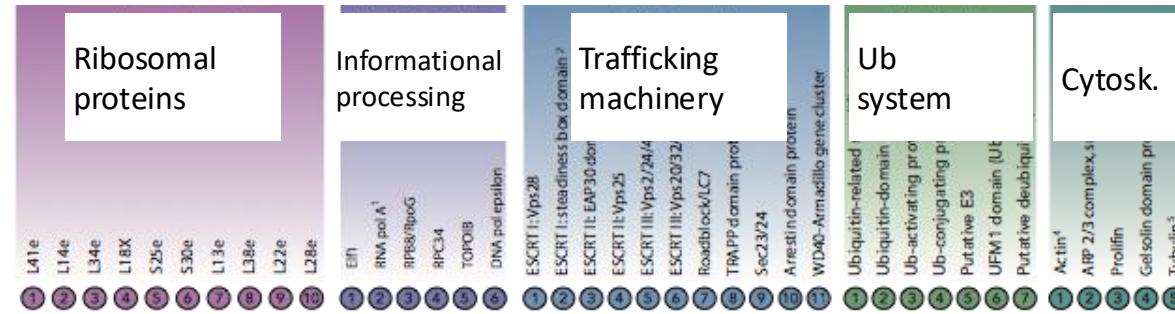
Asgard

Spang 2015, Nature
Zaremba 2017, Nature

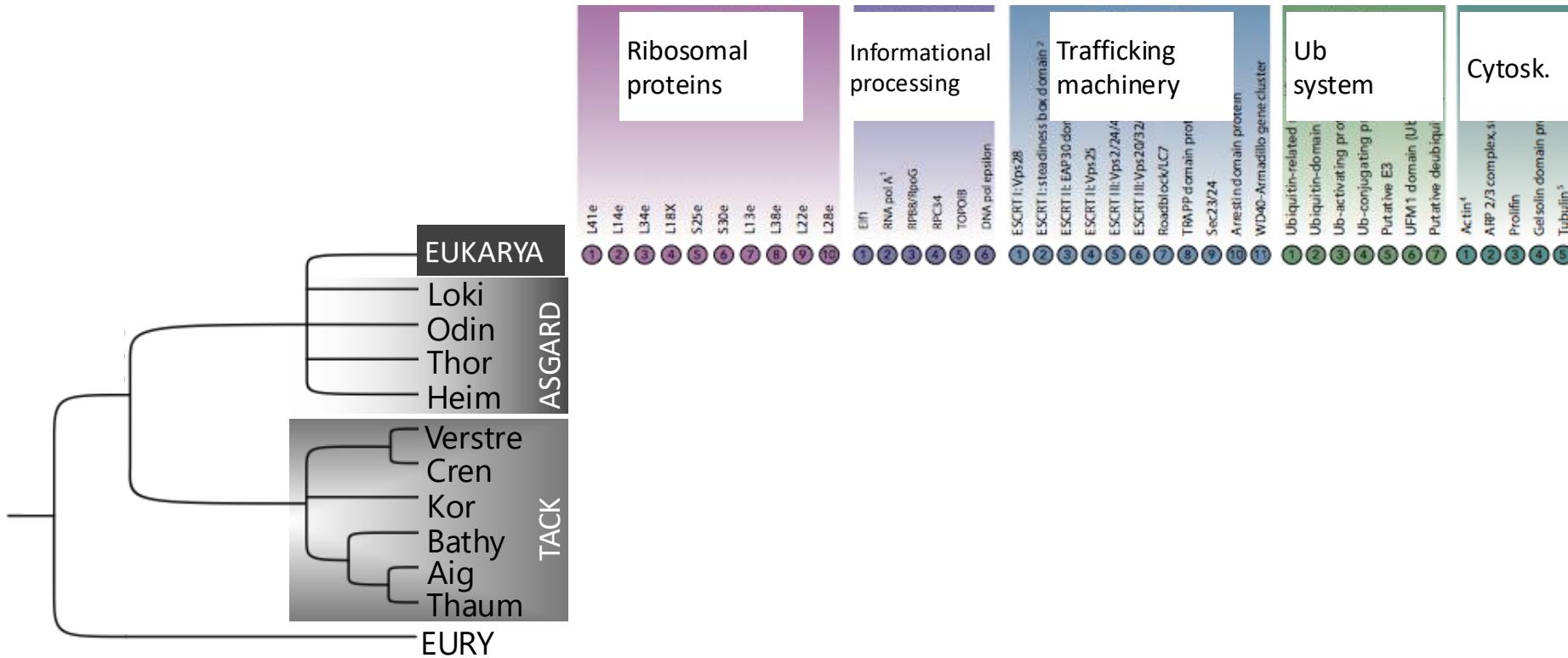
The unveiling of Asgard archaea through metagenomics



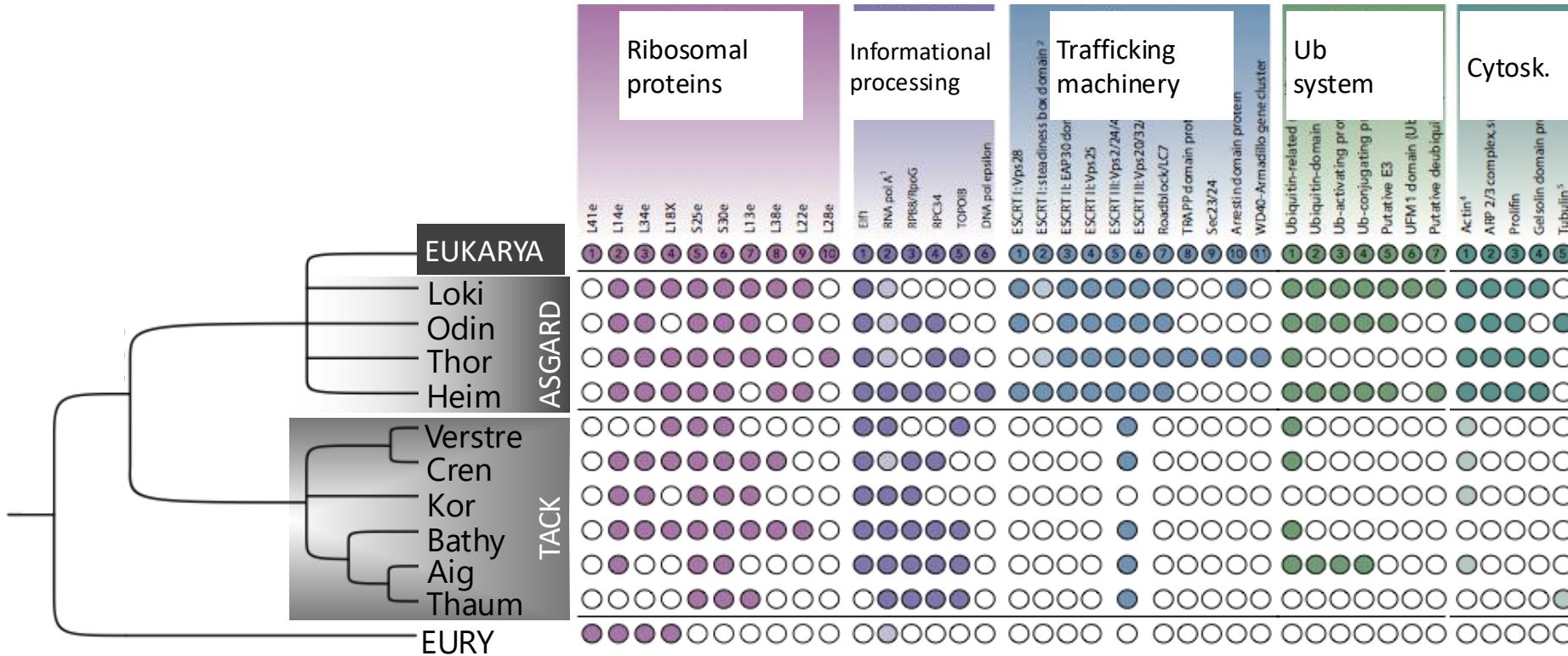
Numerous Eukaryotic Signature Proteins (ESPs) in Asgard archaea



Numerous Eukaryotic Signature Proteins (ESPs) in Asgard archaea



Numerous Eukaryotic Signature Proteins (ESPs) in Asgard archaea



Article

Inference and reconstruction of the heimdallarchaeal ancestry of eukaryotes

<https://doi.org/10.1038/s41586-023-06186-2>

Received: 23 April 2021

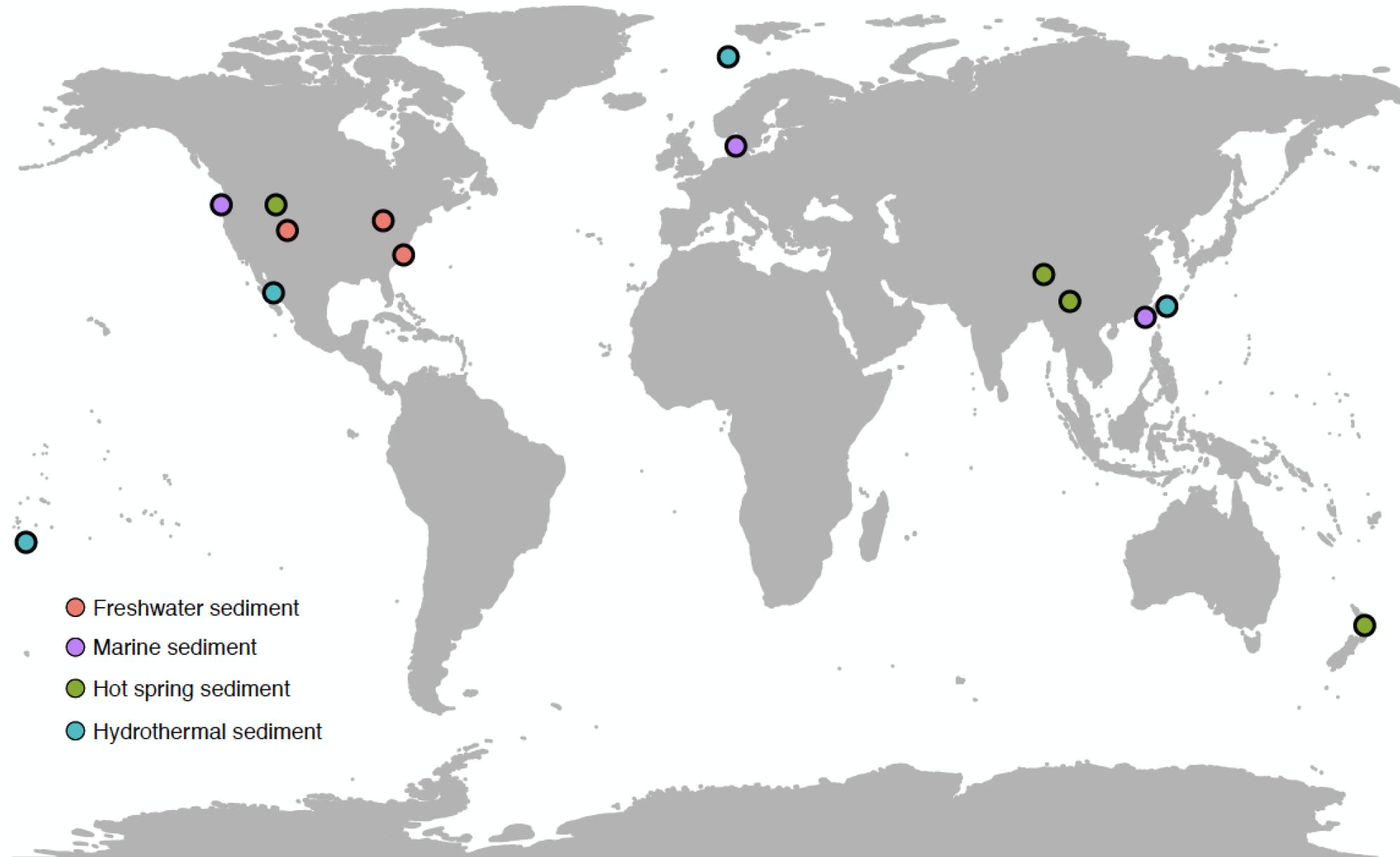
Accepted: 10 May 2023



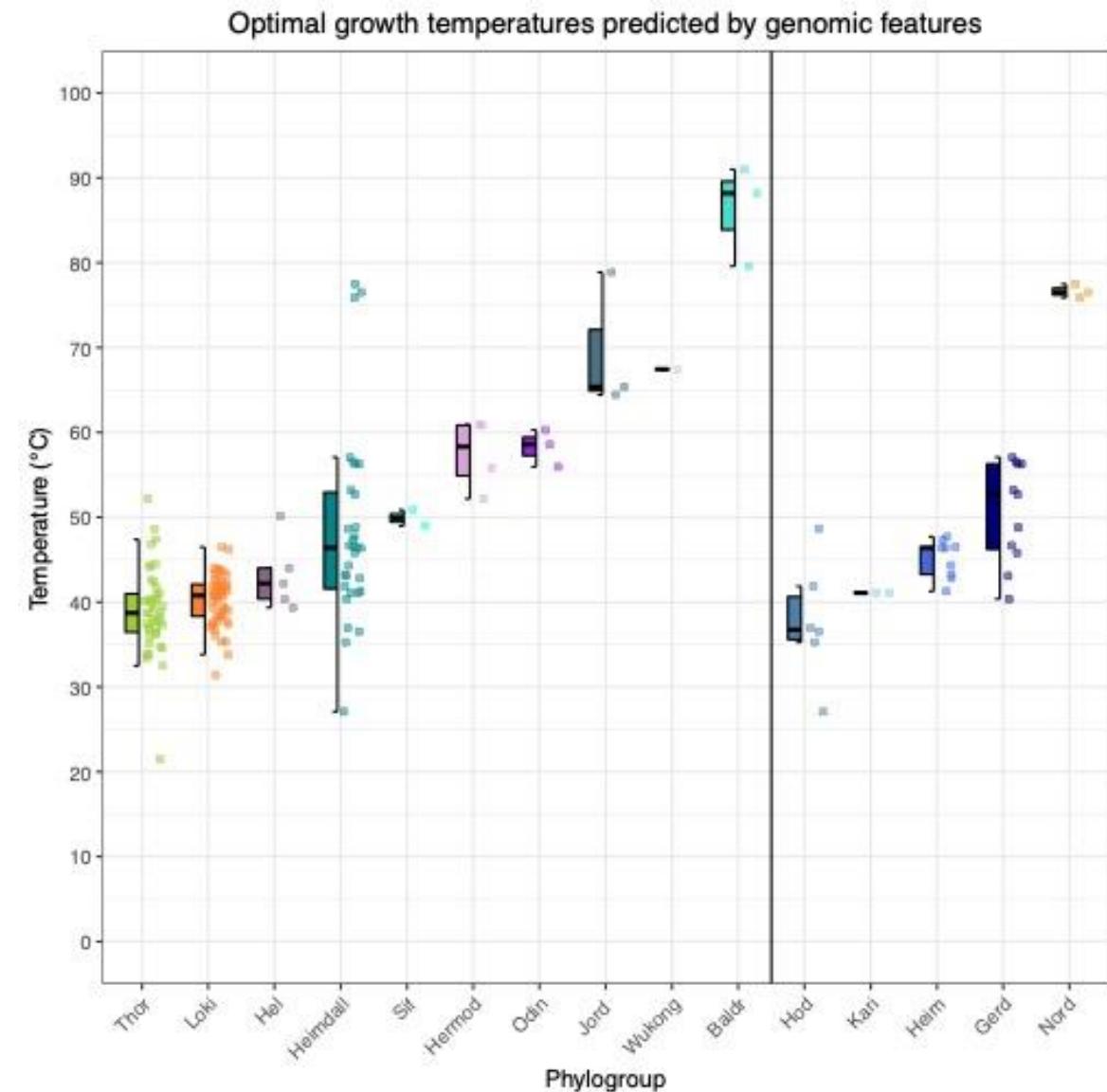
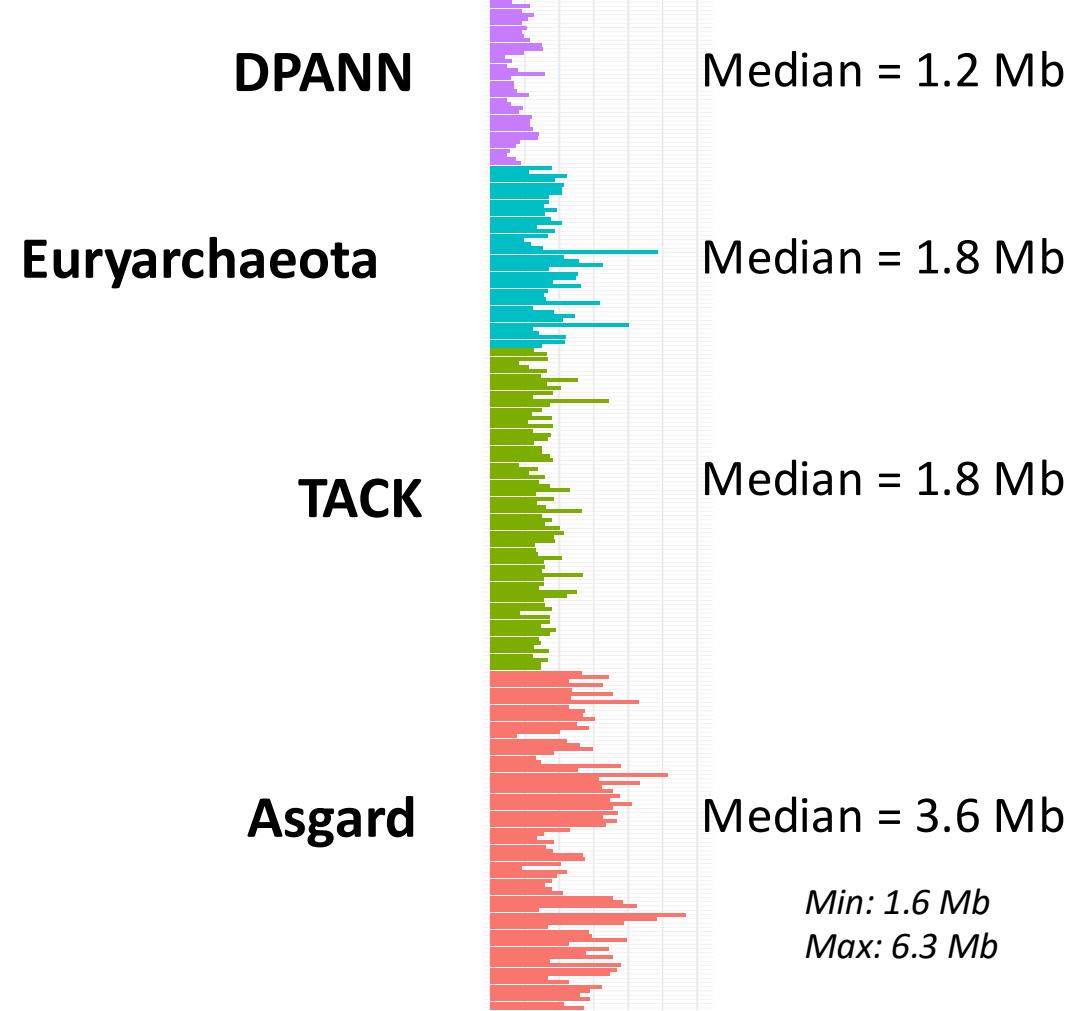
Open access

Laura Eme^{1,2,21}, Daniel Tamarit^{1,3,4,15,21}, Eva F. Caceres^{1,3,21}, Courtney W. Stairs^{1,16}, Valerie De Anda⁵, Max E. Schön¹, Kiley W. Seitz^{5,17}, Nina Dombrowski^{5,18}, William H. Lewis^{1,3,19}, Felix Homa³, Jimmy H. Saw^{1,20}, Jonathan Lombard¹, Takuro Nunoura⁶, Wen-Jun Li⁷, Zheng-Shuang Hua⁸, Lin-Xing Chen⁹, Jillian F. Banfield^{9,10}, Emily St John¹¹, Anna-Louise Reysenbach¹¹, Matthew B. Stott¹², Andreas Schramm¹³, Kasper U. Kjeldsen¹³, Andreas P. Teske¹⁴, Brett J. Baker⁵ & Thijs J. G. Ettema^{1,3}✉

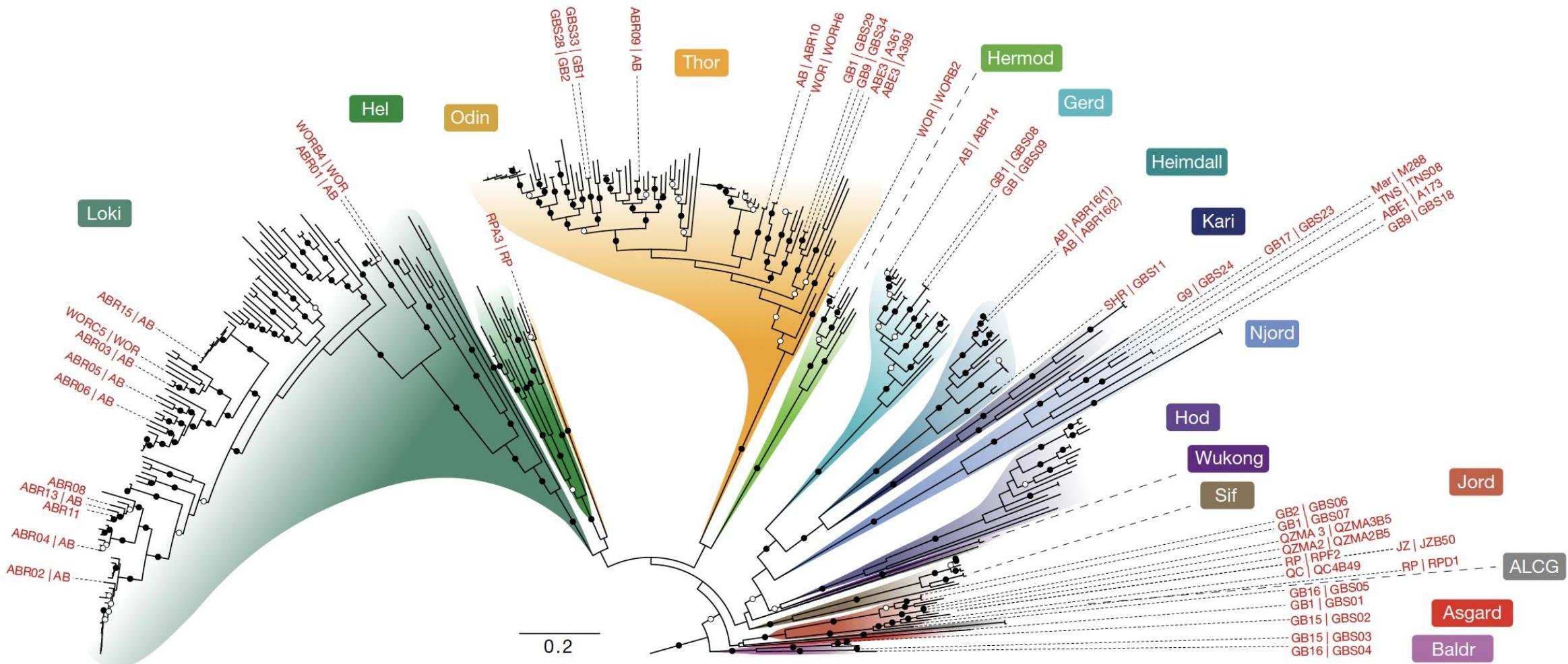
69 new Asgard genomes



Asgard genomes are substantially larger than most archaea



69 new Asgard MAGs

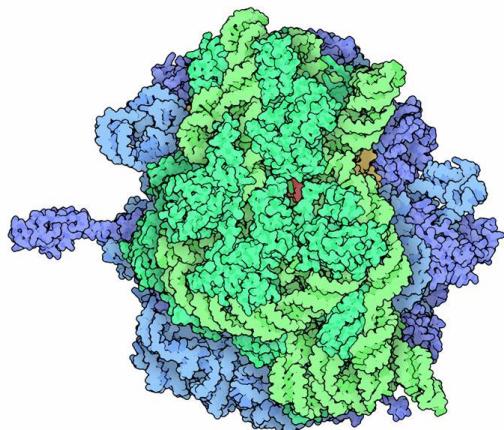


How do Eukaryotes relate to Asgards?

Ribosomal proteins:

- Slow evolving
- Universal
- Coevolve
- Short

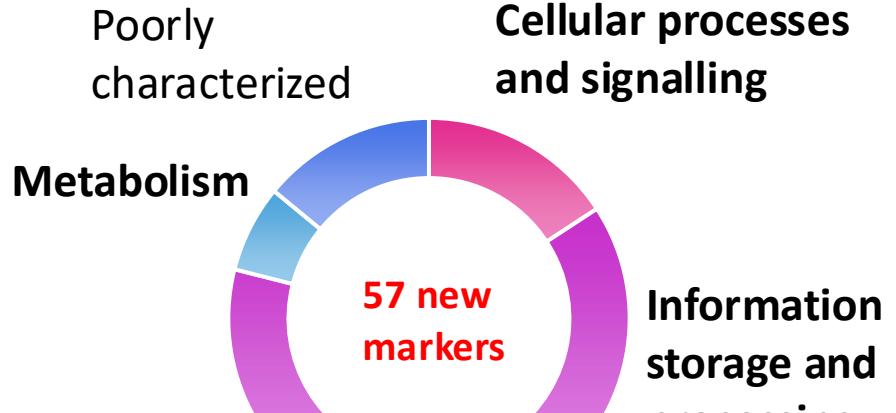
54 ribosomal proteins



~6000 aa

New markers:

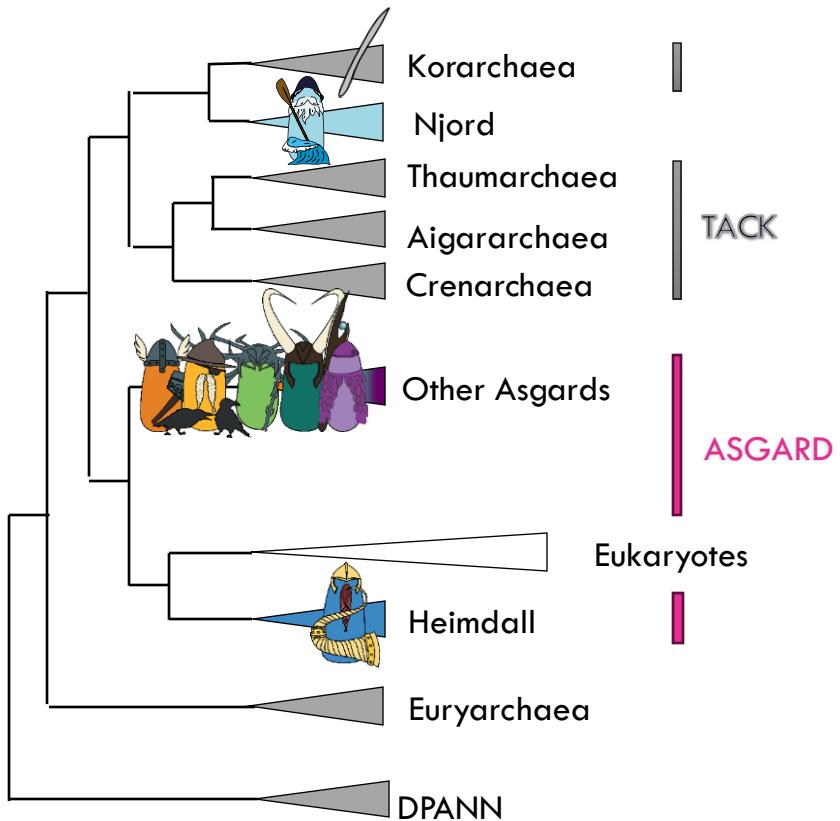
- Conserved in eukaryote and archaea
- Not transferred horizontally



~14,000 aa

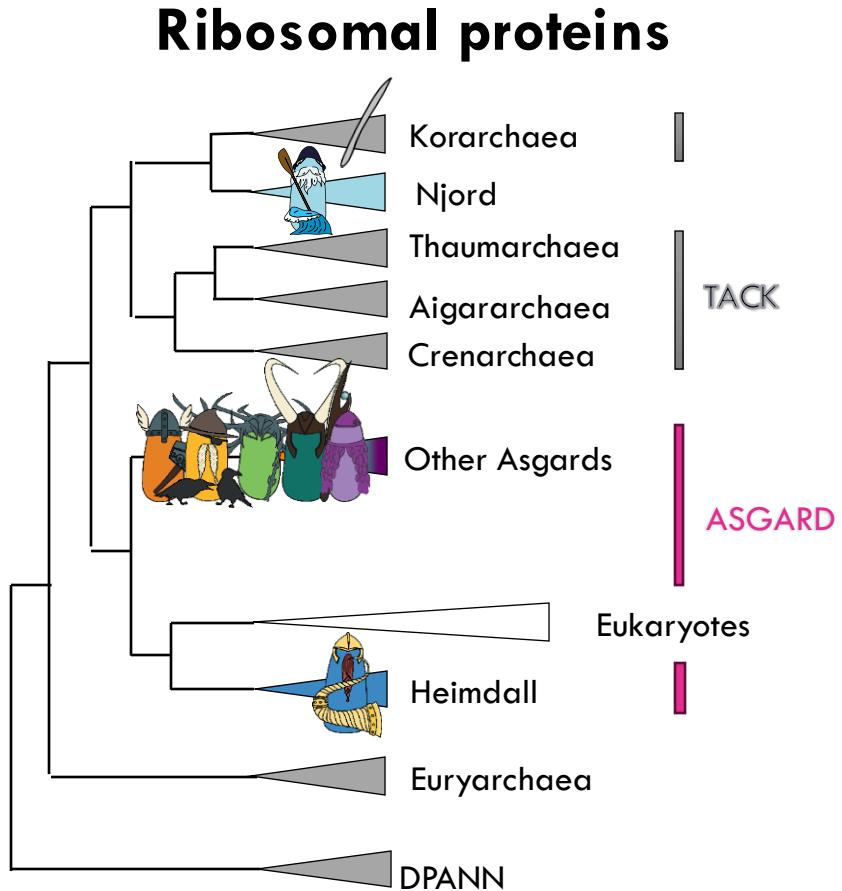
How do Eukaryotes relate to Asgards?

Ribosomal proteins



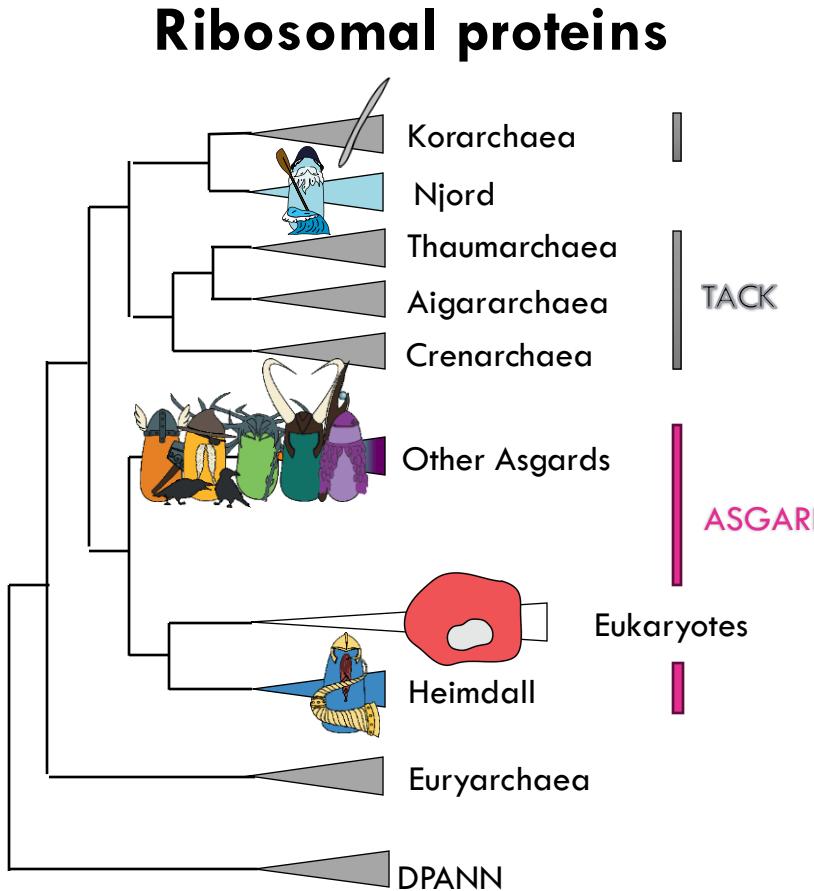
How do Eukaryotes relate to Asgards?

Many ESPs:
Actin,
RPL28e, ...

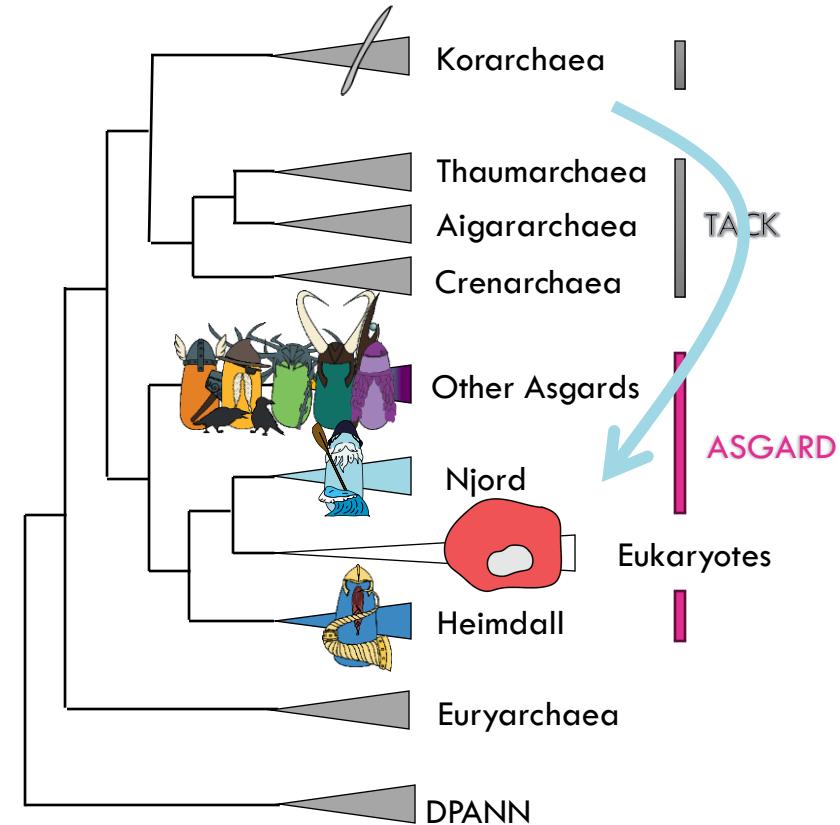


How do Eukaryotes relate to Asgards?

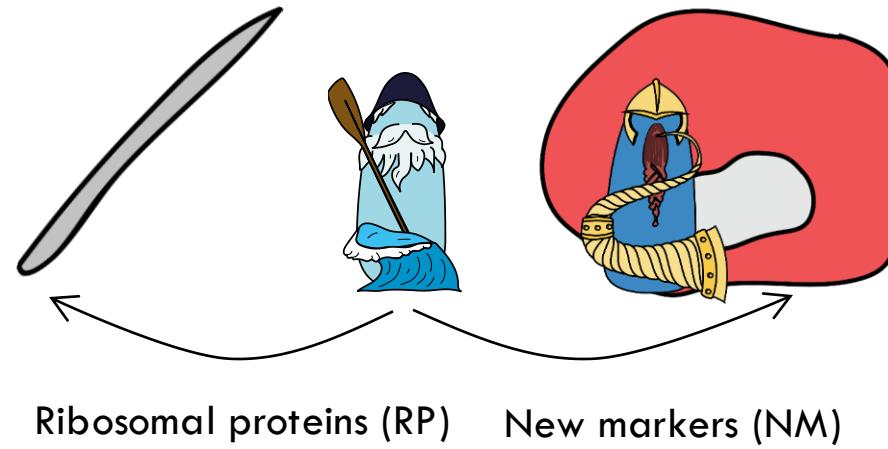
Many ESPs:
Actin,
RPL28e, ...



New markers



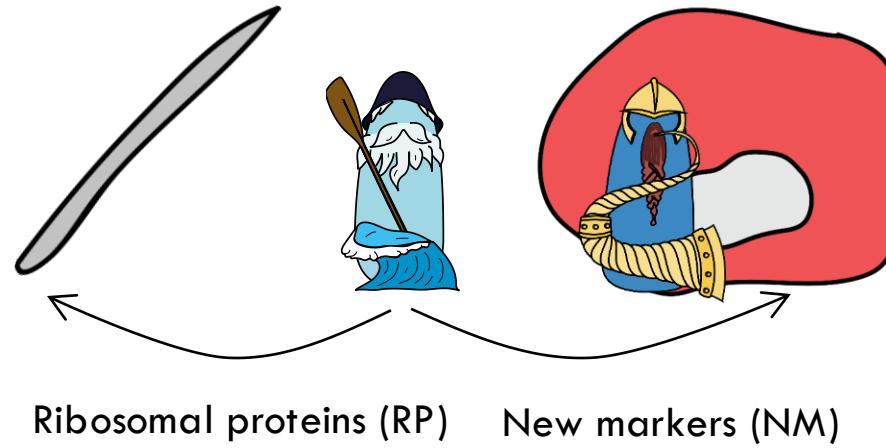
How do Eukaryotes relate to Asgards?



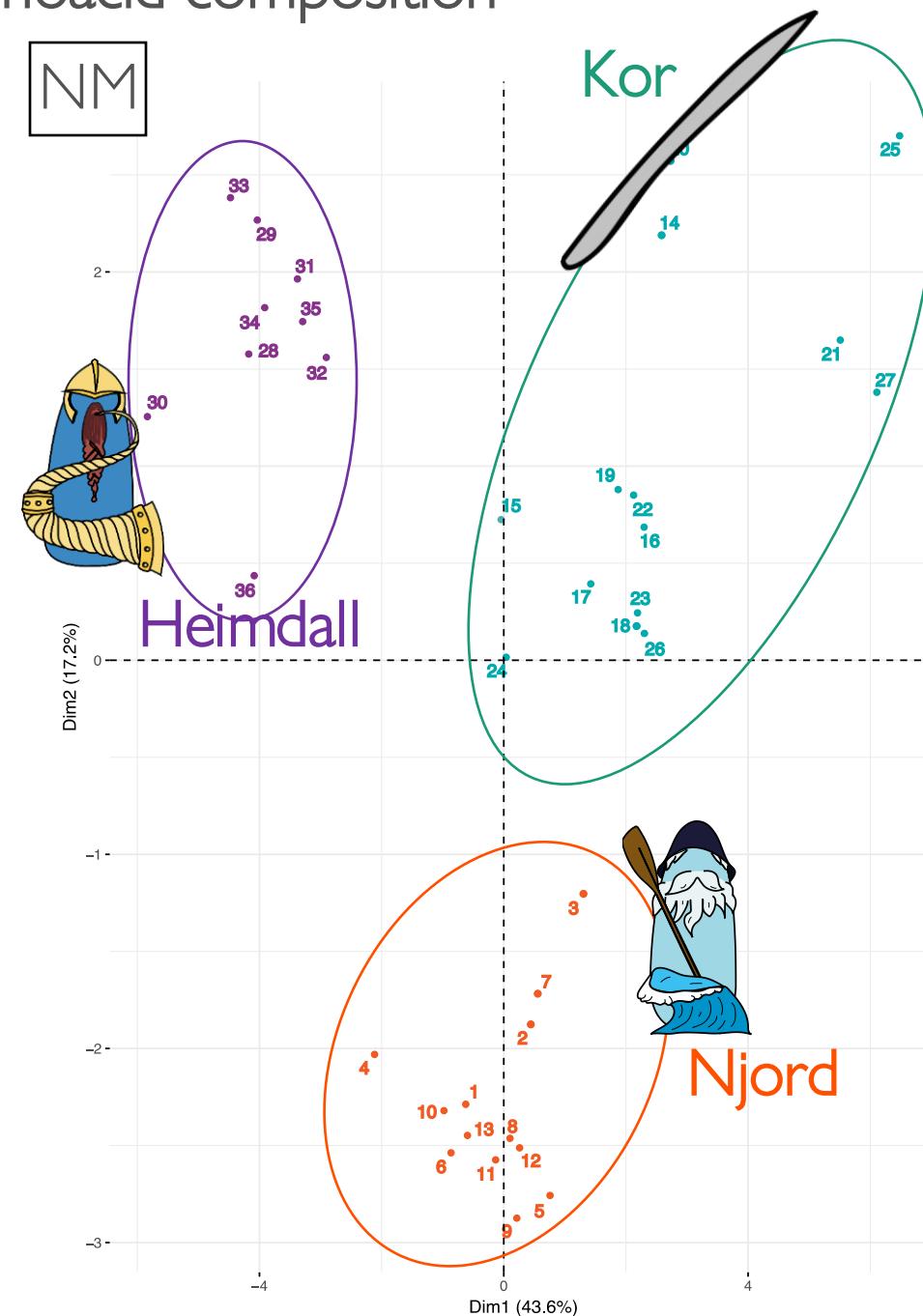
How do Eukaryotes relate to Asgards?

Dataset
Ribosomal proteins
New markers
Software
IQ-Tree (LG+G4+C60+F+PMSF)
Phylobayes (CAT+LG+G4)
Taxon sampling
+/- DPANN
+/- Eukaryotes
+/- Korarchaea
Fast-evolving site removal (FSR)
Recoding (SR4)

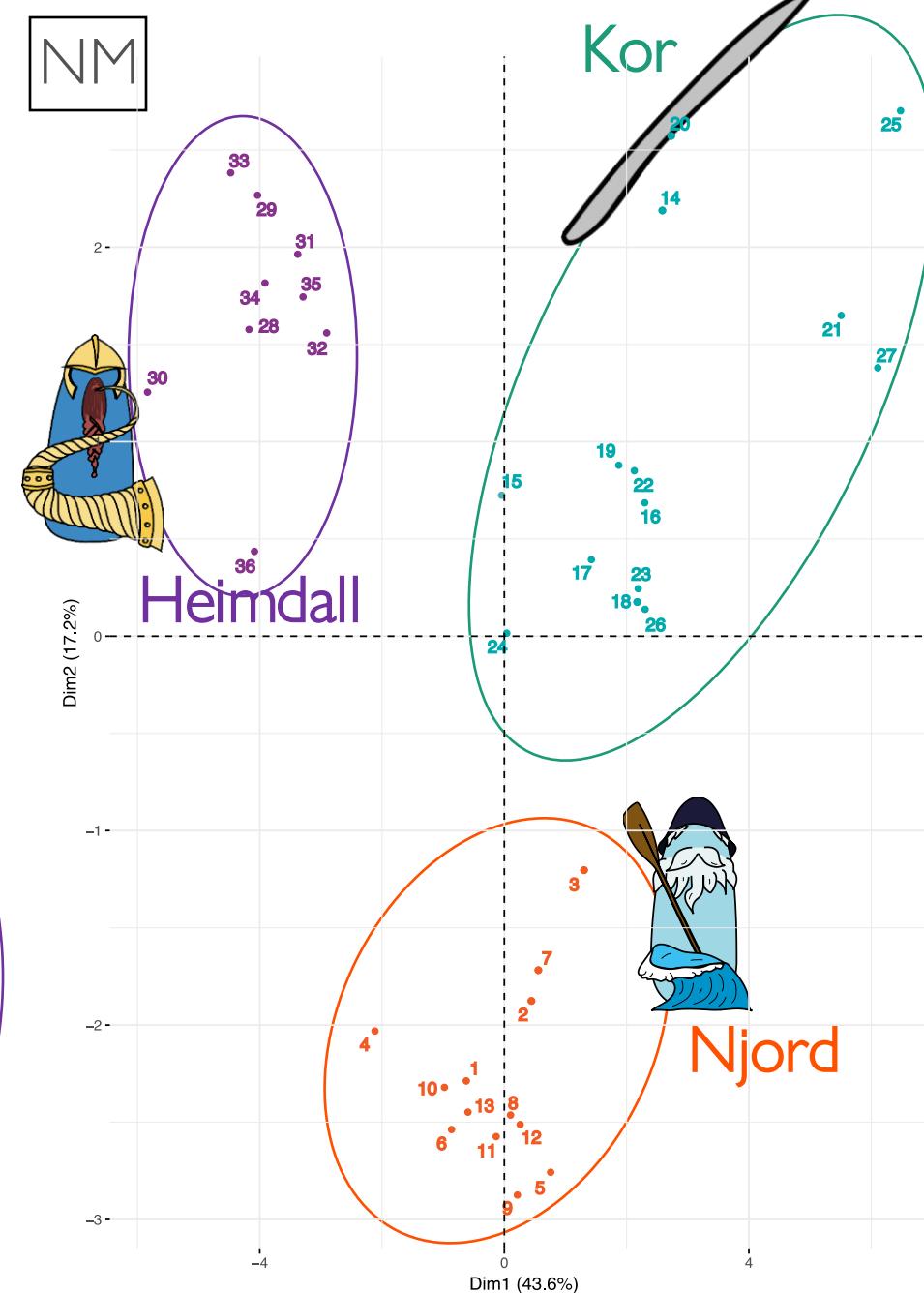
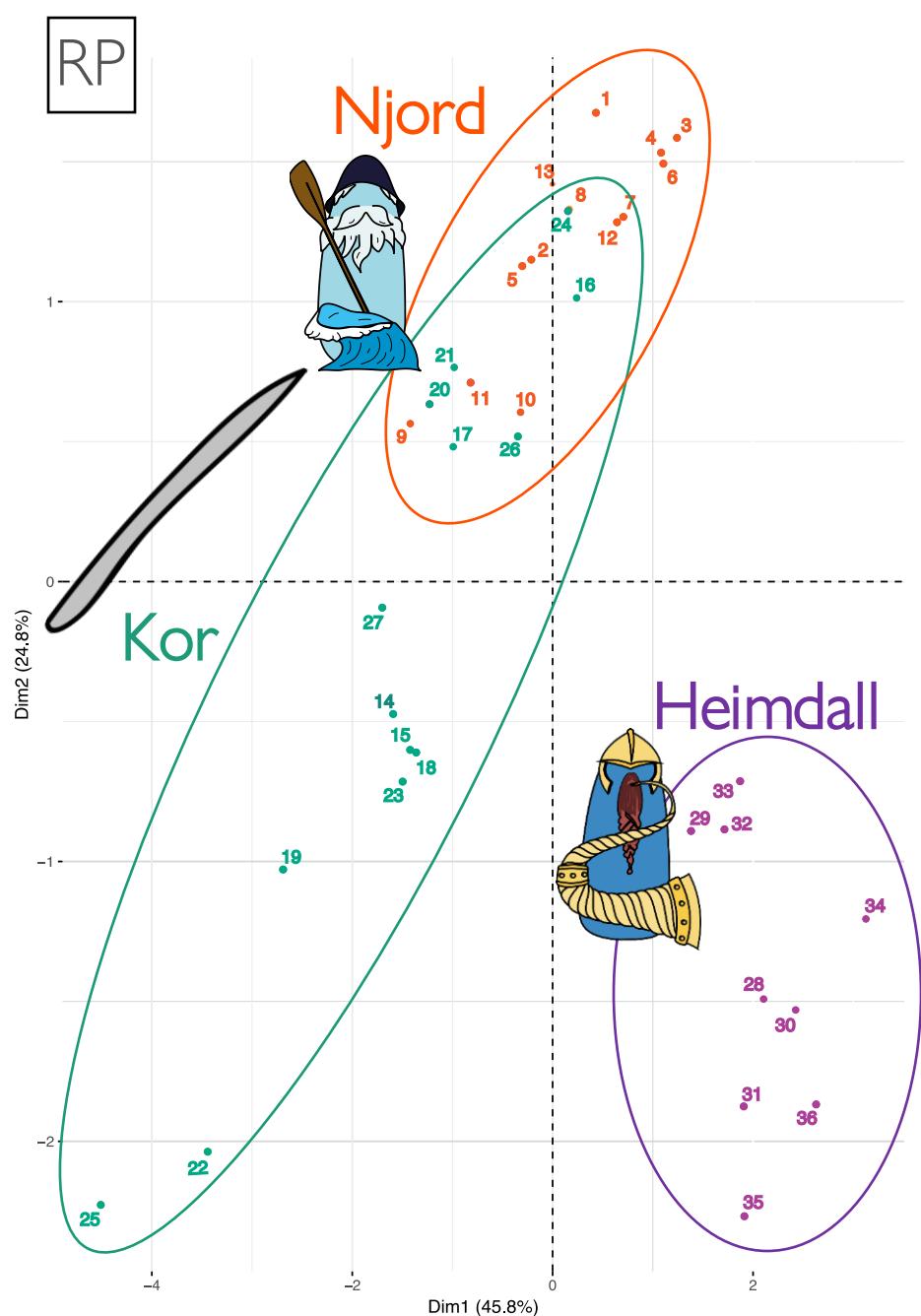
~800 phylogenies...



PCA based on aminoacid composition



PCA based on aminoacid composition



PCA based on aminoacid composition

RP

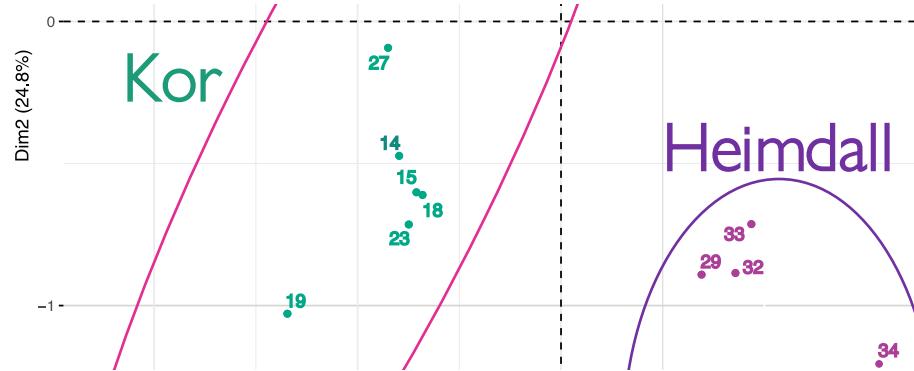
Njord

NM

Kor

Njord and Kor hyperthermophiles:

- rRNA composition bias
- structural constraints in the ribosome explaining convergent AA composition



NM

Heimdall

Kor

25

21

27

14

19

22

16

17

23

18

26

24

36

3

-2

-1

0

1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36

3

1

4

6

8

10

12

13

15

17

19

21

23

25

27

29

31

33

35

37

39

41

43

45

47

49

51

53

55

57

59

61

63

65

67

69

71

73

75

77

79

81

83

85

87

89

91

93

95

97

99

101

103

105

107

109

111

113

115

117

119

121

123

125

127

129

131

133

135

137

139

141

143

145

147

149

151

153

155

157

159

161

163

165

167

169

171

173

175

177

179

181

183

185

187

189

191

193

195

197

199

201

203

205

207

209

211

213

215

217

219

221

223

225

227

229

231

233

235

237

239

241

243

245

247

249

251

253

255

257

259

261

263

265

267

269

271

273

275

277

279

281

283

285

287

289

291

293

295

297

299

301

303

305

307

309

311

313

315

317

319

321

323

325

327

329

331

333

335

337

339

341

343

345

347

349

351

353

355

357

359

361

363

365

367

369

371

373

375

377

379

381

383

385

387

389

391

393

395

397

399

401

403

405

407

409

411

413

415

417

419

421

423

425

427

429

431

433

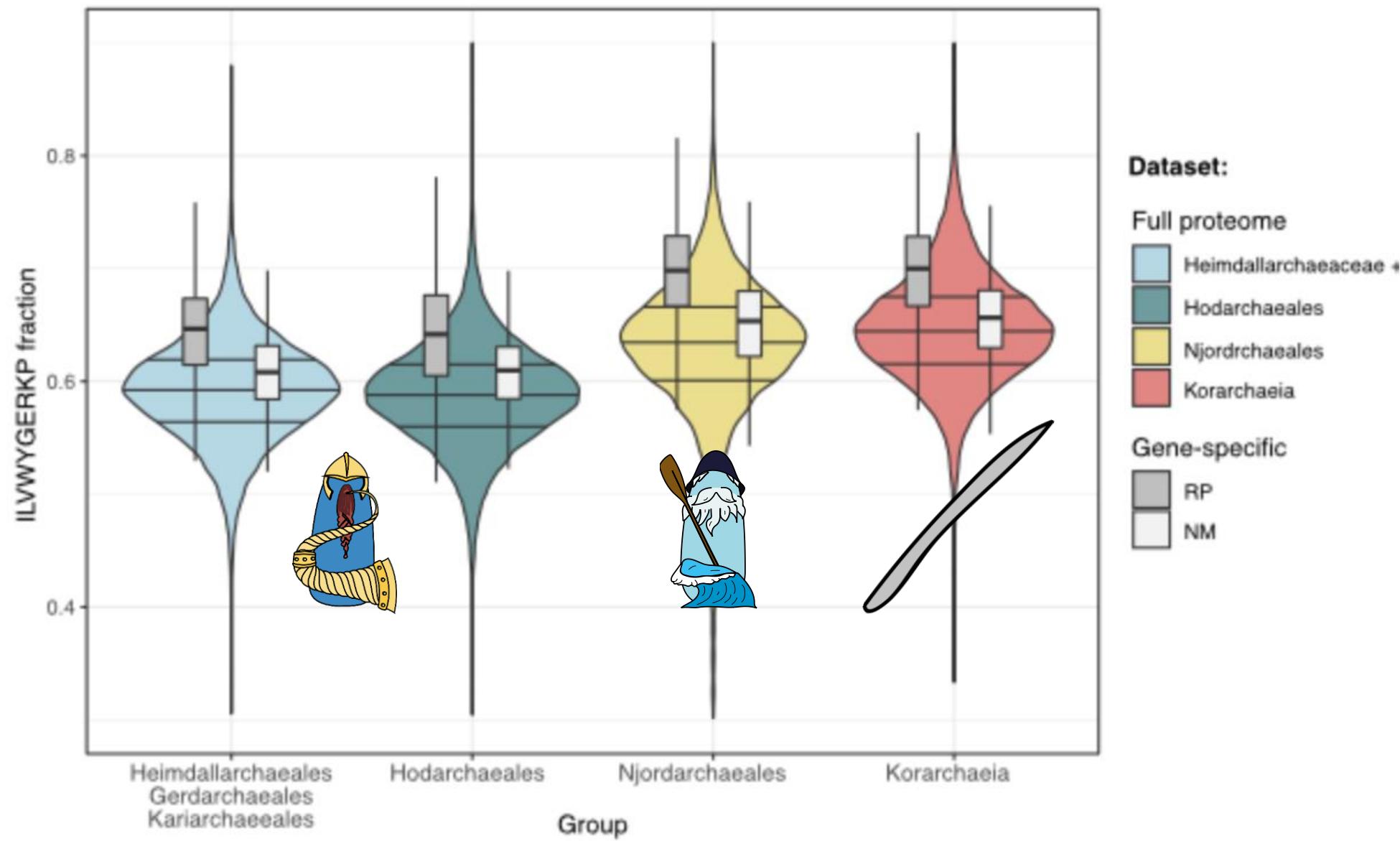
435

437

439

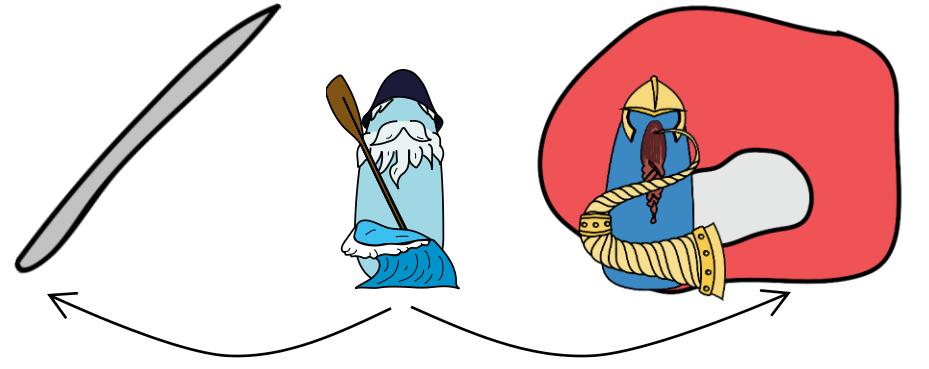
B

Compositional bias: ILVWYGERKP fraction



How do Eukaryotes relate to Asgards?

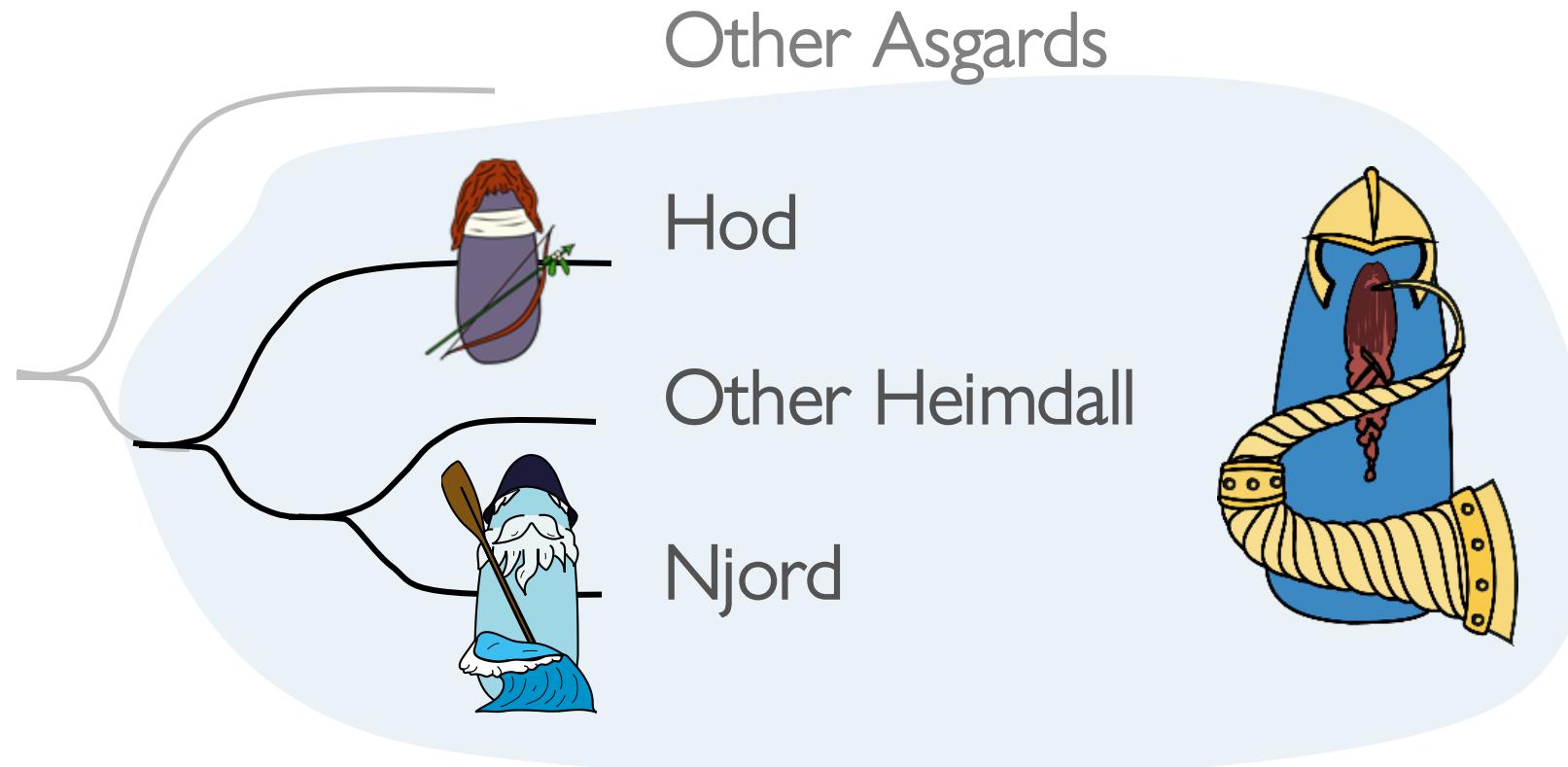
Dataset	Ribosomal proteins New markers
Software	IQ-Tree (LG+G4+C60+F+PMSF) Phylobayes (CAT+LG+G4)
Taxon sampling	+/- DPANN +/- Eukaryotes +/- Korarchaea
Fast-evolving site removal (FSR)	
Recoding (SR4)	



Ribosomal proteins (RP) New markers (NM)
RP (without Korarchaea)

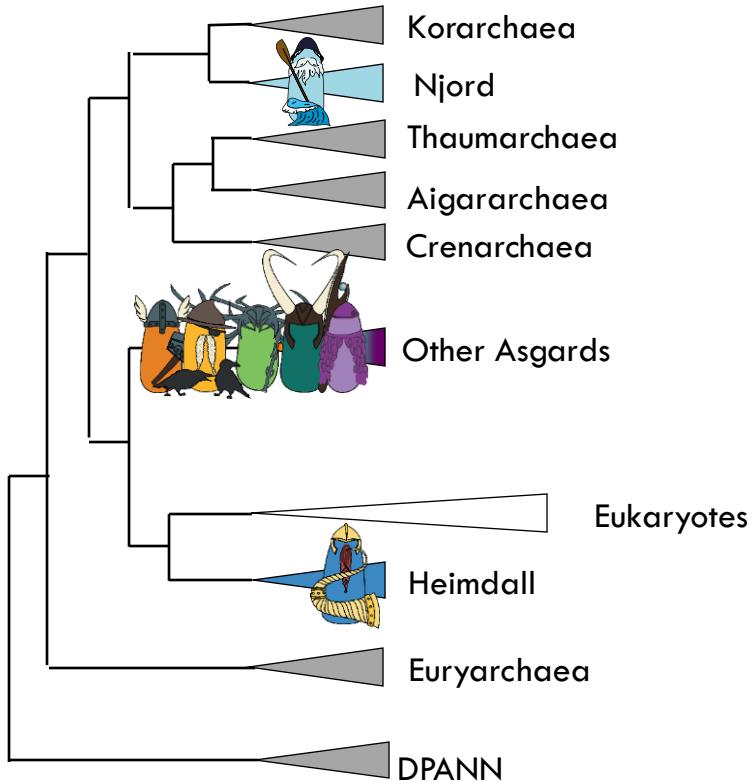
~800 phylogenies...

Njord are in fact a sub-clade of Heimdallarchaeia

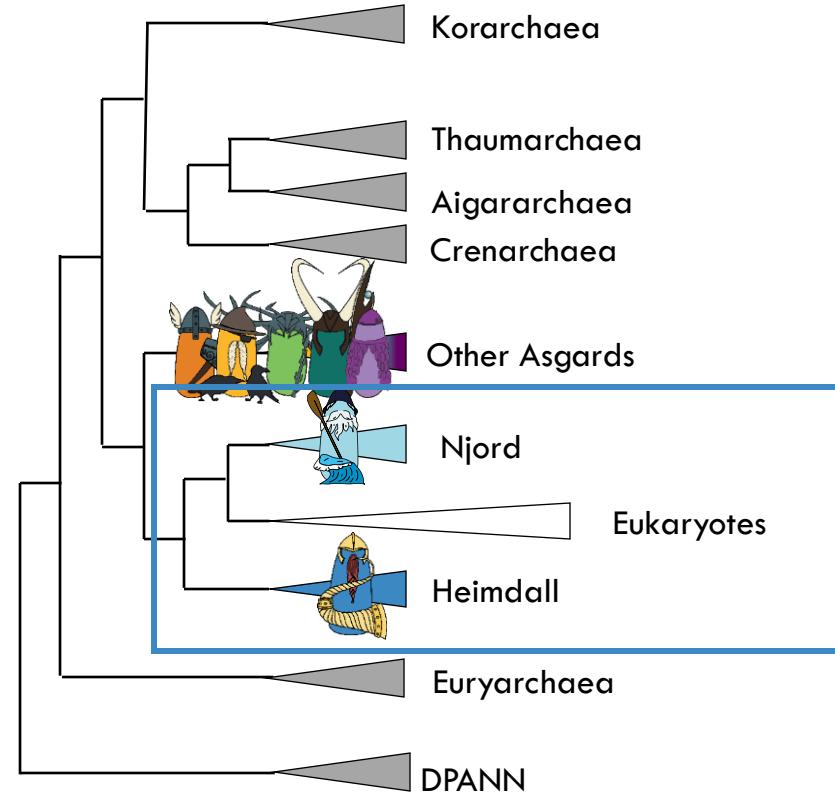


Why do we care, again?

Ribosomal proteins



New markers

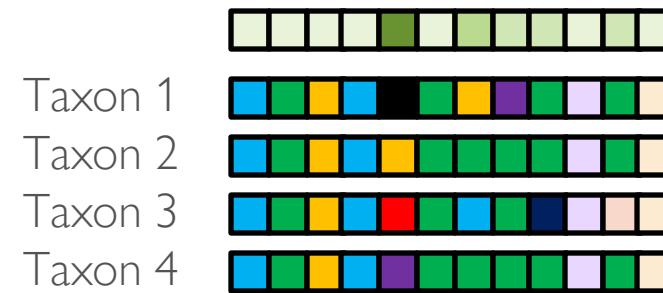
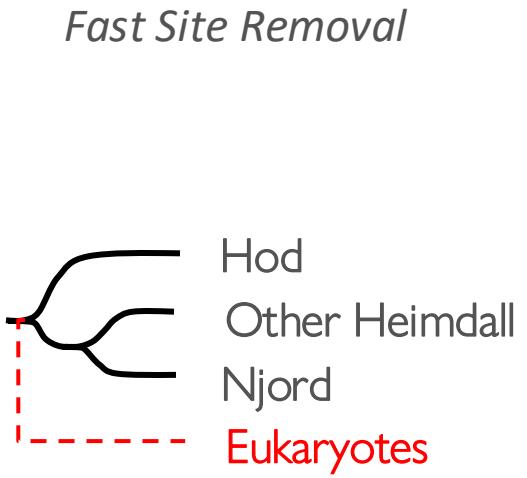
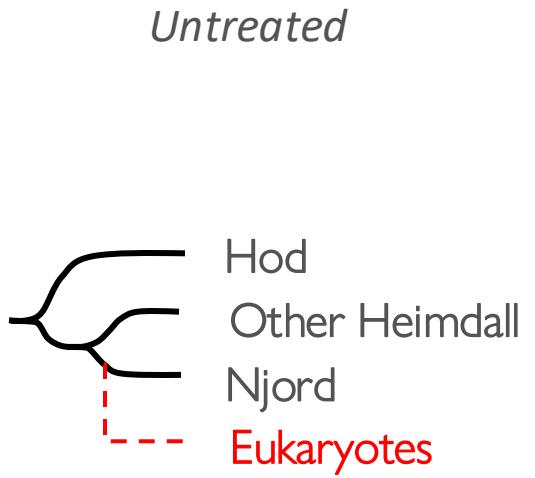


Placing Eukaryotes in the Asgard tree require careful phylogenetic investigations

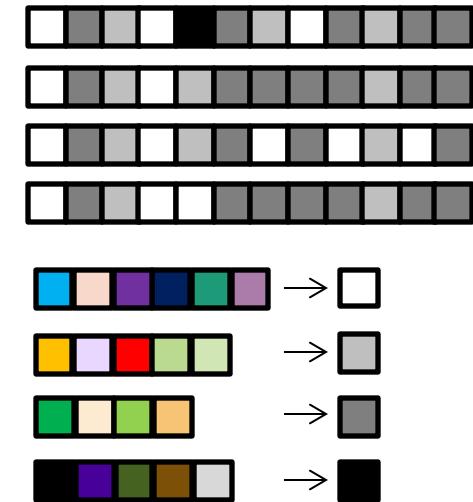
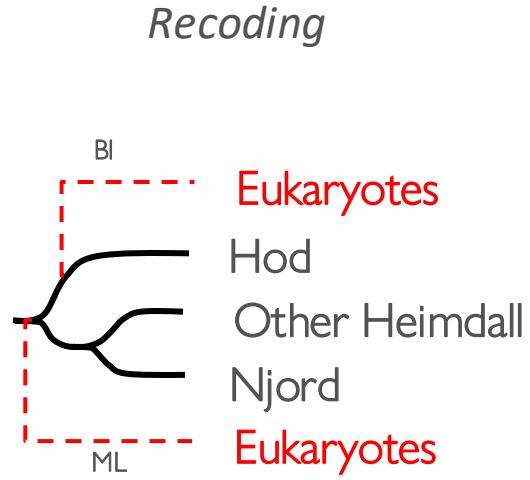
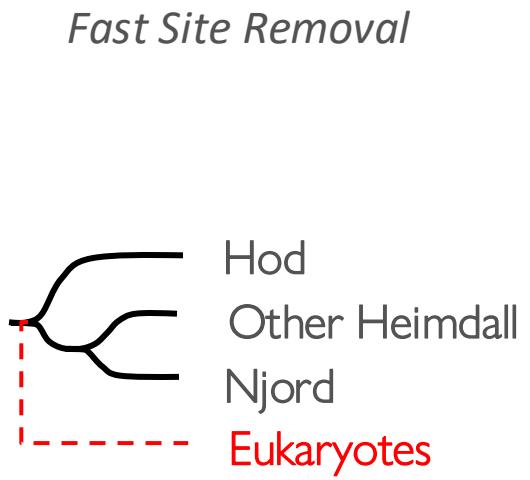
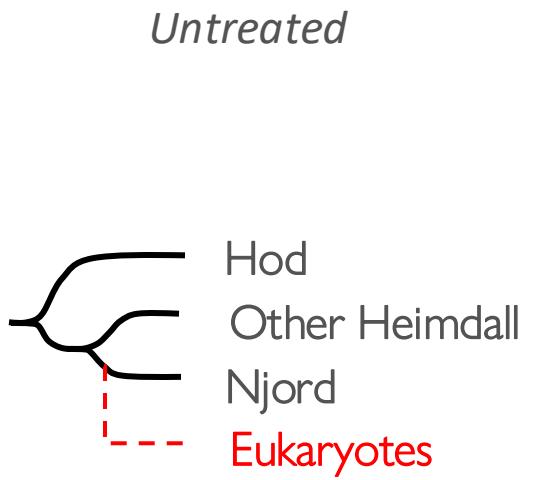
Untreated



Placing Eukaryotes in the Asgard tree require careful phylogenetic investigations

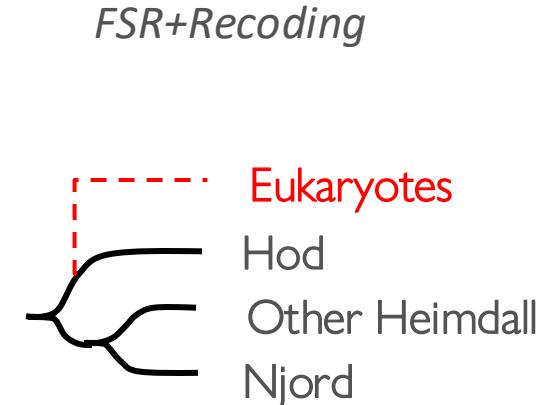
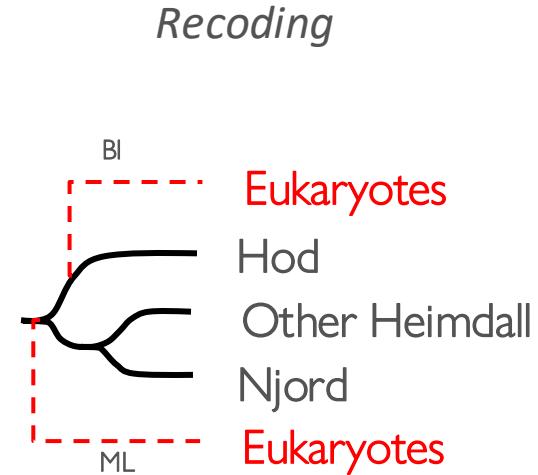
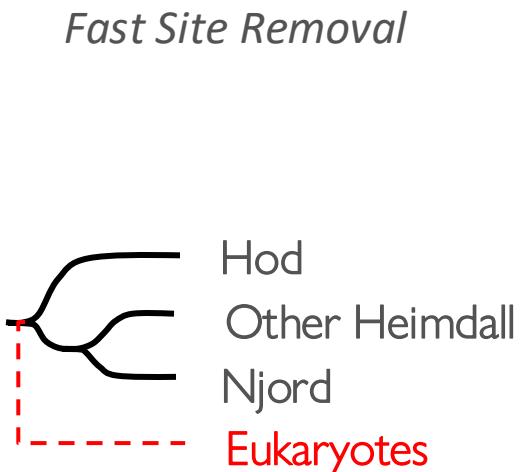
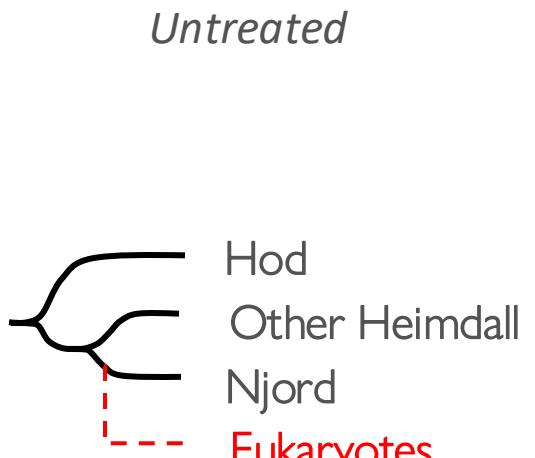


Placing Eukaryotes in the Asgard tree require careful phylogenetic investigations

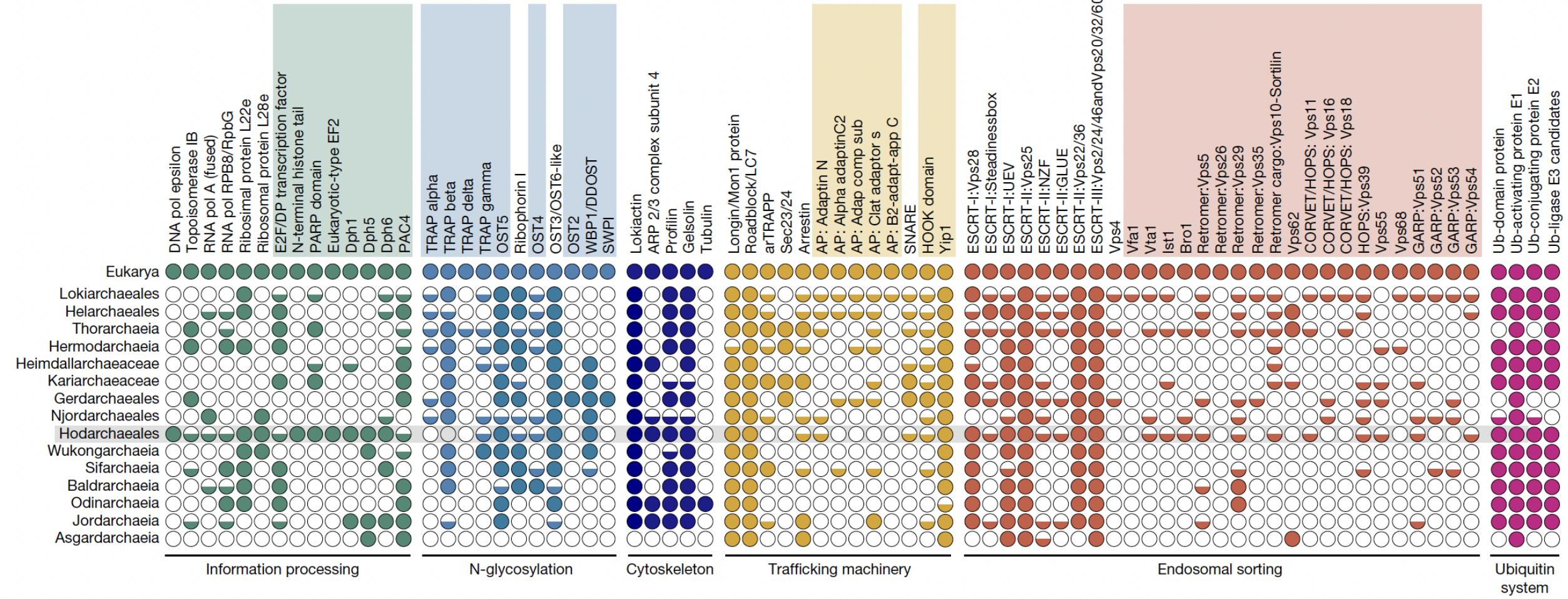


E.g.: SR4 recoding (Susko and Roger 2007)

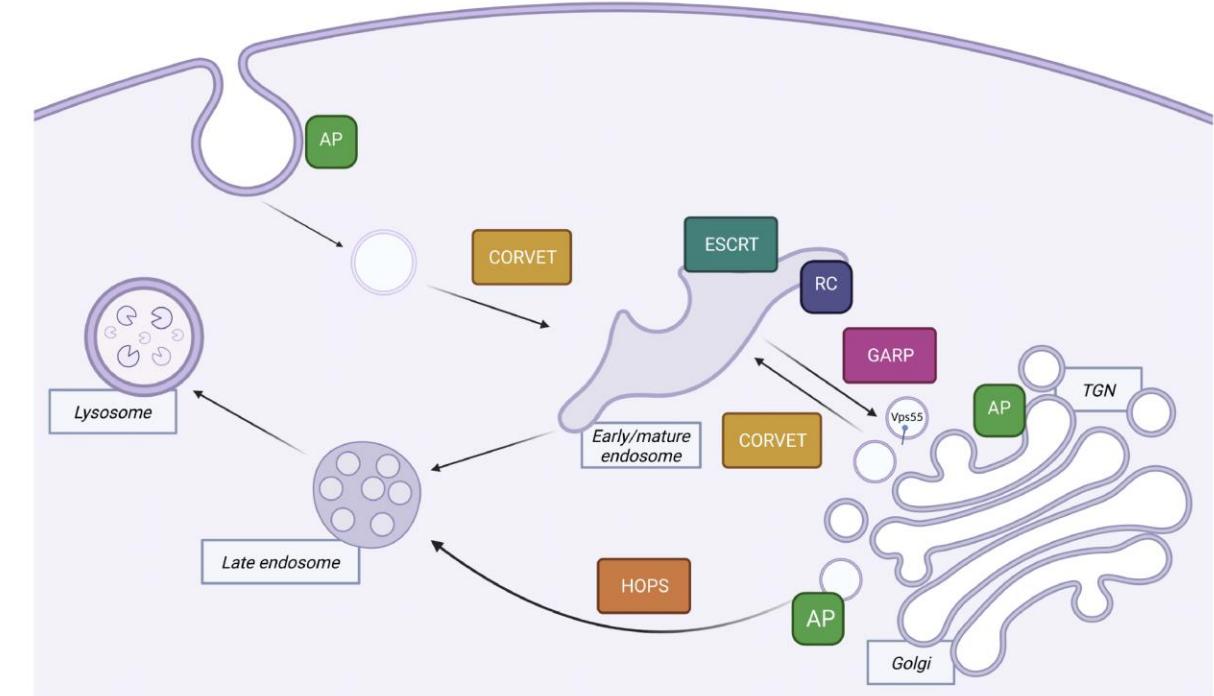
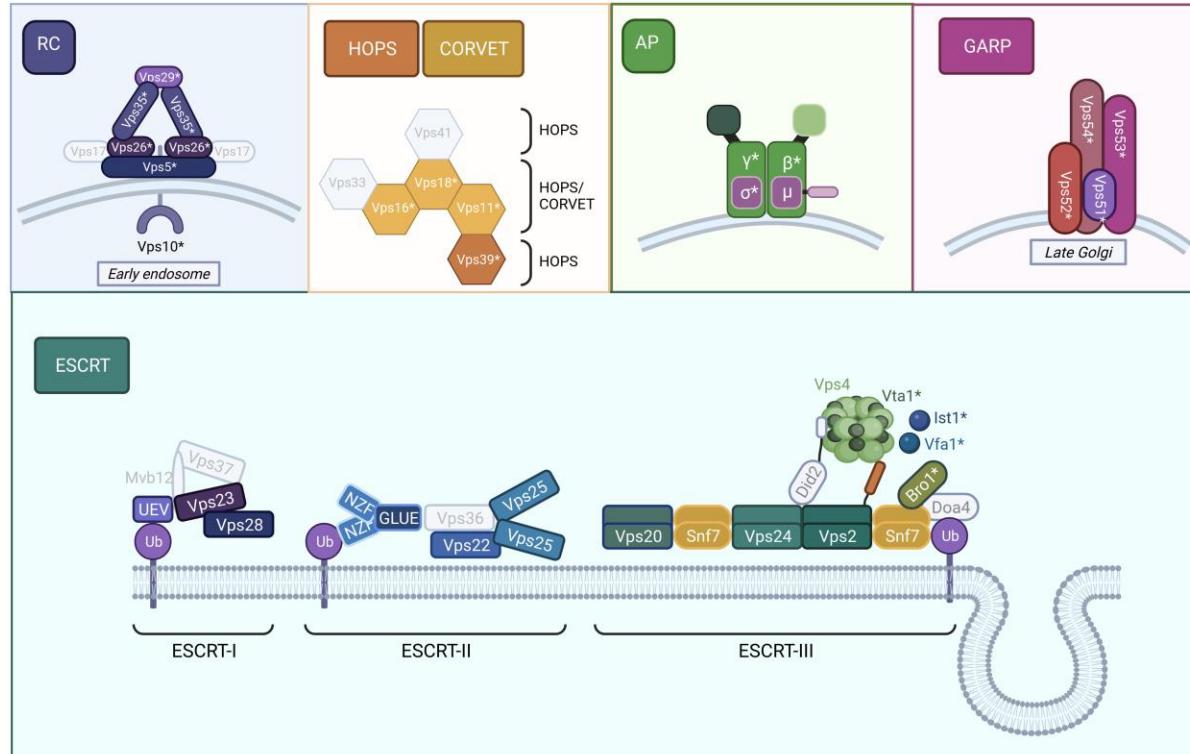
Placing Eukaryotes in the Asgard tree require careful phylogenetic investigations



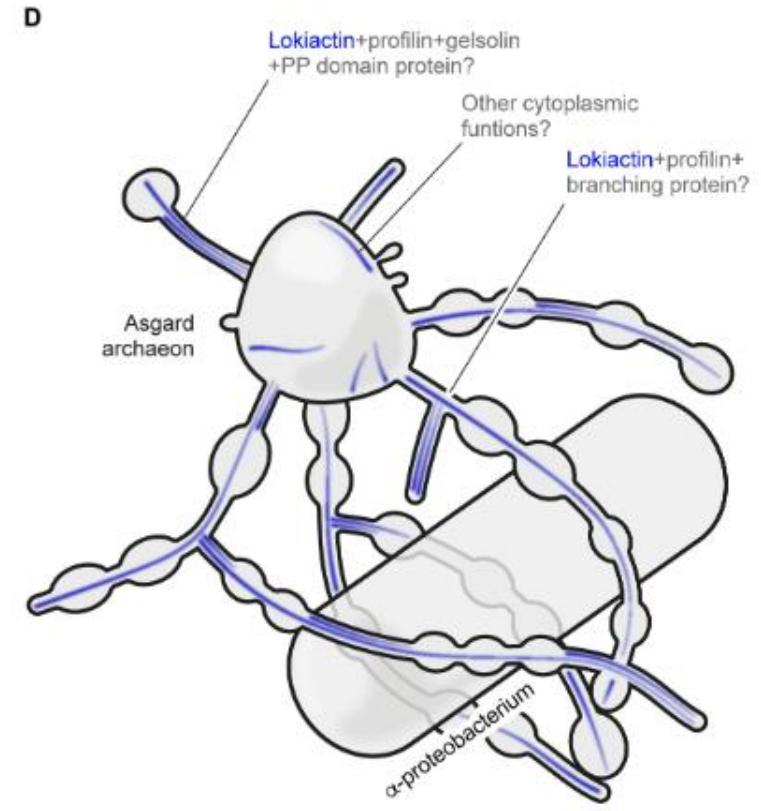
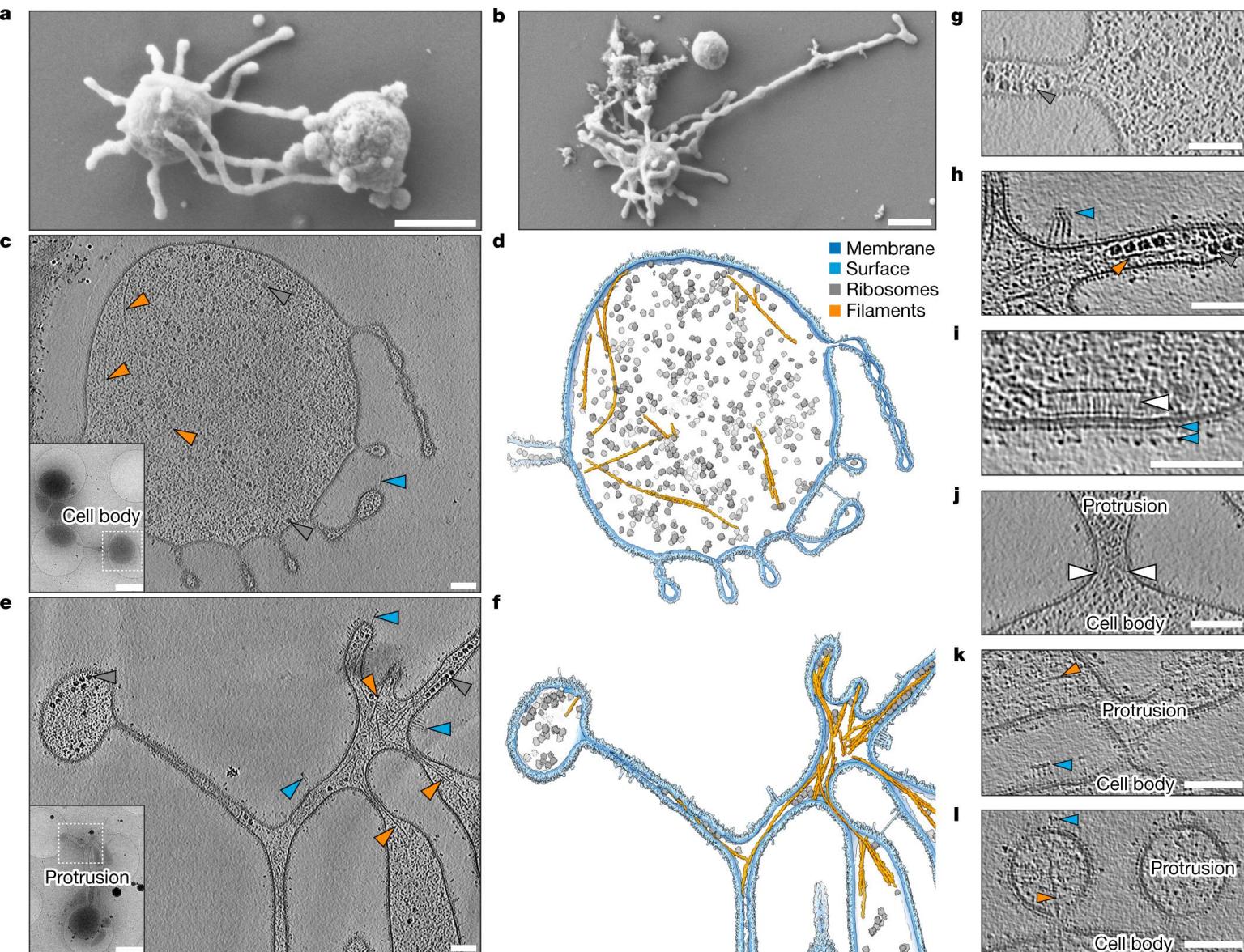
New ESPs support the relationship of Eukaryotes and Hod



New ESPs involved in complex intracellular trafficking



Asgard archaea have an actin-based cytoskeleton



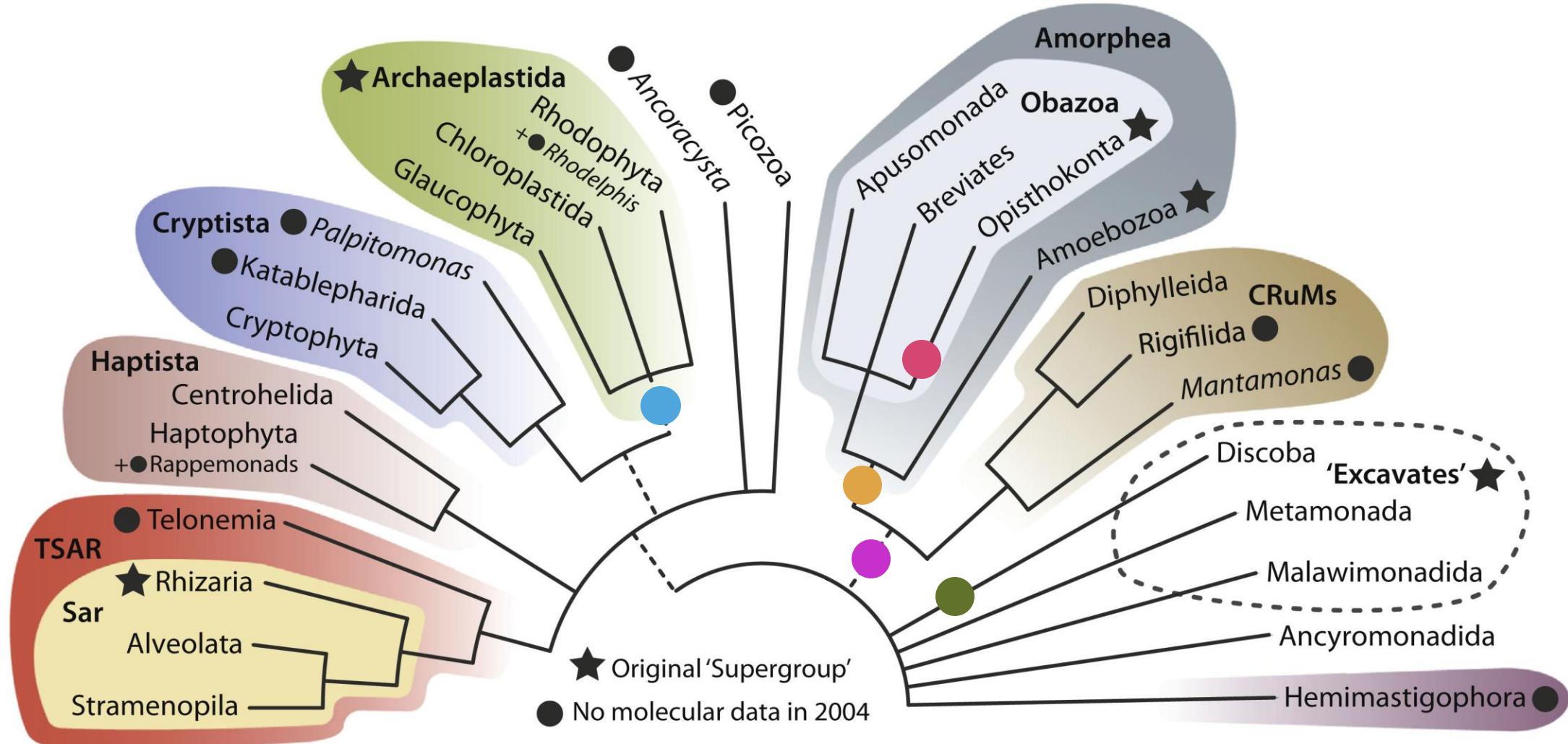
Charles-Orszag A, Petek-Seoane NA, Mullins RD. 2024. Archaeal actins and the origin of a multi-functional cytoskeleton. *J Bacteriol*

Rodrigues-Oliveira, T., Wollweber, F., Ponce-Toledo, R.I. et al. Actin cytoskeleton and complex cell architecture in an Asgard archaeon. *Nature* **613**, 332–339 (2023).

Example 3

Rooting the tree of eukaryotes

Williamson, Eme, et al. *Nature*, 2025



Trends in Ecology & Evolution

Modified from Burki et al., 2020

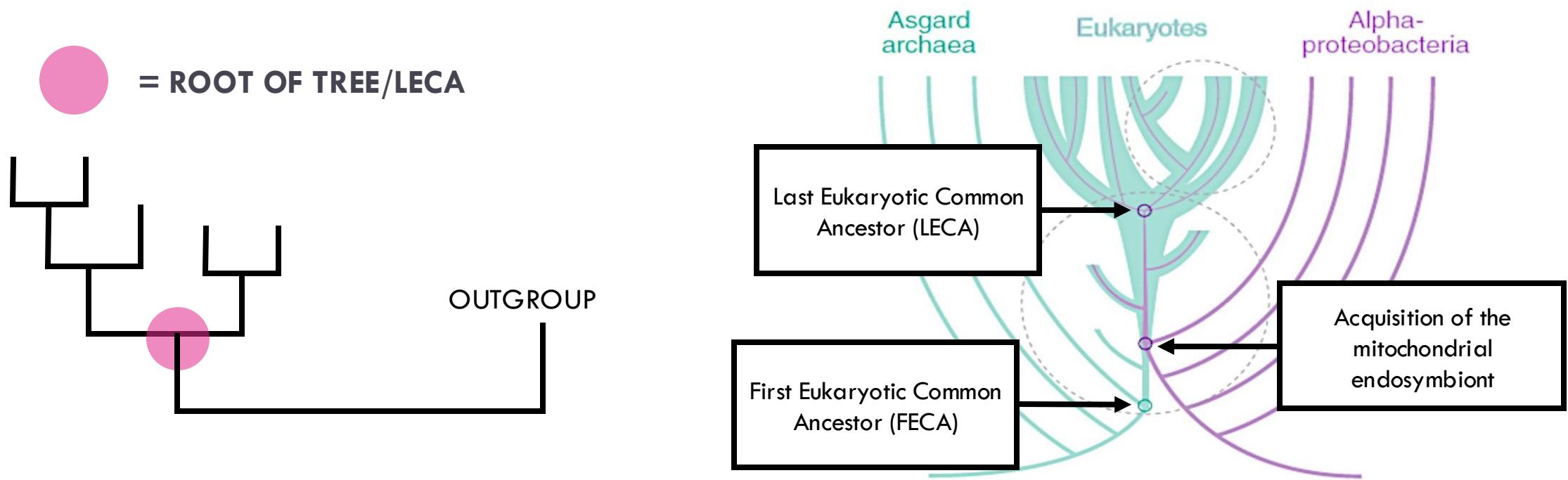
Gene tree parsimony, gene fusions and domain evolution, molecular and cellular features, replacement of highly conserved amino acids...

PROBLEMS:

Convergence/reversion, lack of evolutionary models, no tests
for robustness, no consensus,
sensitive to missing data and incomplete taxonomic sampling

OUTGROUP ROOTING

Outgroups are lineages that fall outside of the group being studied, but are closely related enough to retain phylogenetic signal through orthologous genes



The evolutionary relationship of eukaryotes with archaea and alphaproteobacteria make them possible outgroups

Shorter branch length between alpha and eukaryotes

Modified from Roger et al., 2017

Identify proteins of alphaproteobacterial origin in a small set of eukaryotes

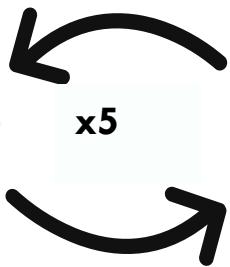


Retrieve homologs from selected eukaryotes and prokaryotes



Kelsey Williamson

Manual removal of non-orthologs and outliers



Estimate single gene trees for each marker gene (IQ-TREE)



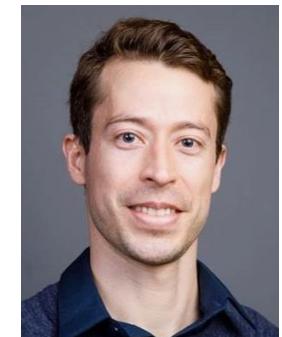
Laura Eme

FINAL DATASETS

93 marker genes
63 eukaryotes
37 alphaproteobacteria

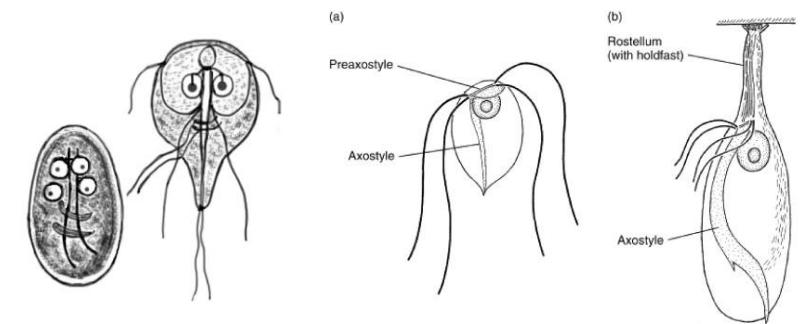


Select outgroup size and composition



Sergio
Muñoz-Gómez

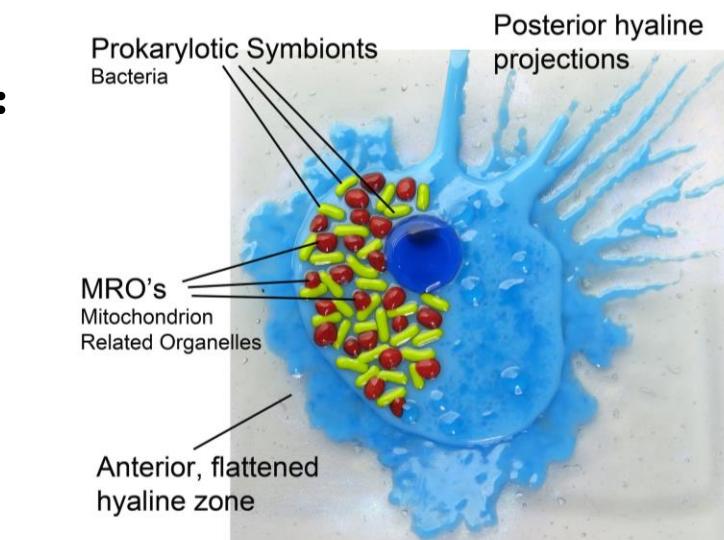
Metamonads are problematic taxa for datasets of mitochondrial proteins



Metamonada consists entirely of **anaerobes**:
lack a mitochondrial genome;
have highly reduced mitochondrion-related organelles
→ lack most of the proteins in the dataset

Metamonads with the most mitochondrial proteins were included:
Anaeramoeba ignava – 13 genes, 19% site occupancy
Anaeramoeba flamelloides – 11 genes, 16.5% site occupancy

Datasets with and without metamonads were generated:
Anae+ and **Anae-**



Anaeramoeba fused glass dish
Jane Hartman

TESTING THE POSITION OF THE ROOT

1. Estimate rooted phylogeny with site-heterogenous models
2. Create alternate root positions to evaluate likelihood differences between previously proposed roots and the optimal root estimated here
3. Evaluate all root positions under additional complex evolutionary models that account for different phenomena

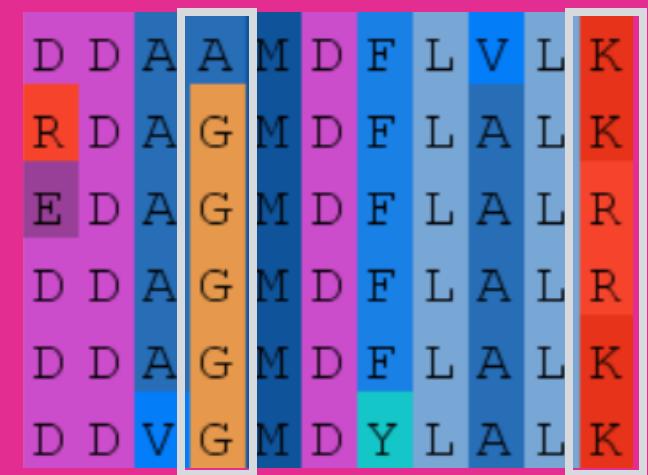
SITE-HETEROGENEOUS MODELS

Site heterogeneous models account for the fact that different sites evolve under different evolutionary constraints.

C-series models (E.g. C60)

UDM models (E.g. UDM64)

Set of amino acid frequency classes generated from public protein databases



SITE-HETEROGENEOUS MODELS

MEOW (MAMMaL Extension On Whole alignment)

- Amino acid frequency classes are estimated directly from the data
- Can use different proportions of high- versus low-rate sites to generate custom site classes

CAT-PMSF

- Generates set of site-specific amino acid frequencies
- Phylobayes is run on a fixed guide tree (ex. LG+G) to estimate frequencies



Hector
Baños



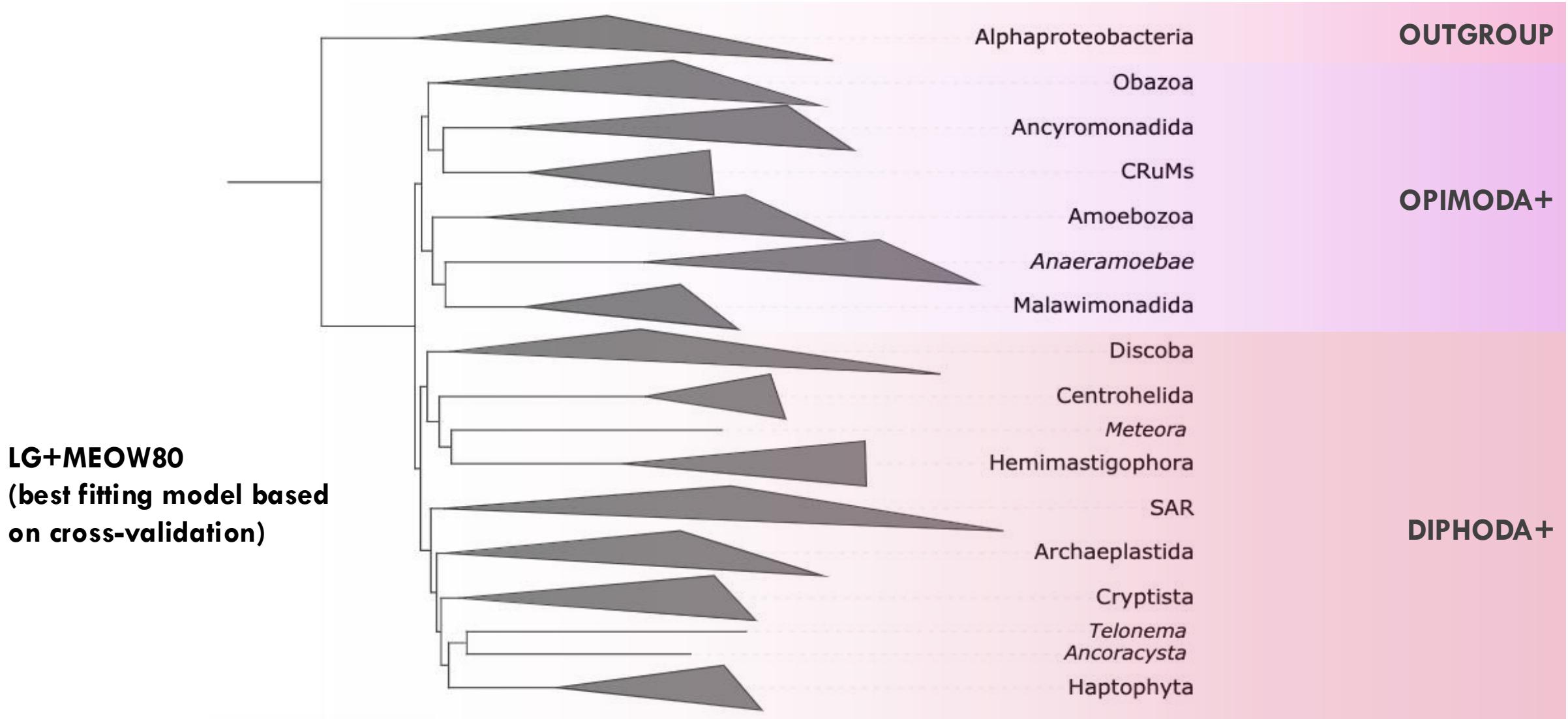
Ed
Susko



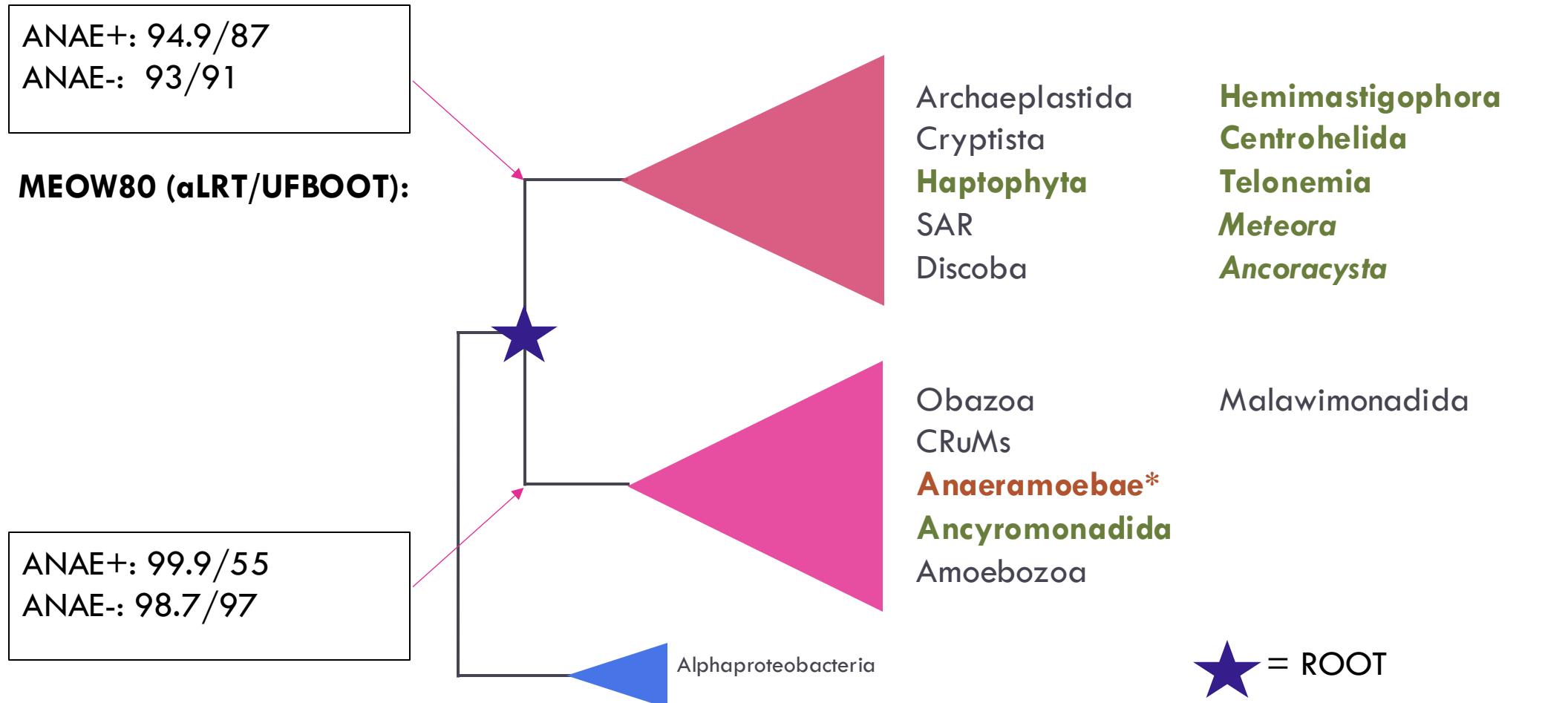
Andrew
Roger

D	D	A	A	M	D	F	L	V	L	K
R	D	A	G	M	D	F	L	A	L	K
E	D	A	G	M	D	F	L	A	L	R
D	D	A	G	M	D	F	L	A	L	R
D	D	A	G	M	D	F	L	A	L	K
D	D	V	G	M	D	Y	L	A	L	K

All models recover a root separating eukaryotes into ‘Opimoda+’ and ‘Diphoda+’



All models recover a root separating eukaryotes into ‘Opimoda+’ and ‘Diphoda+’

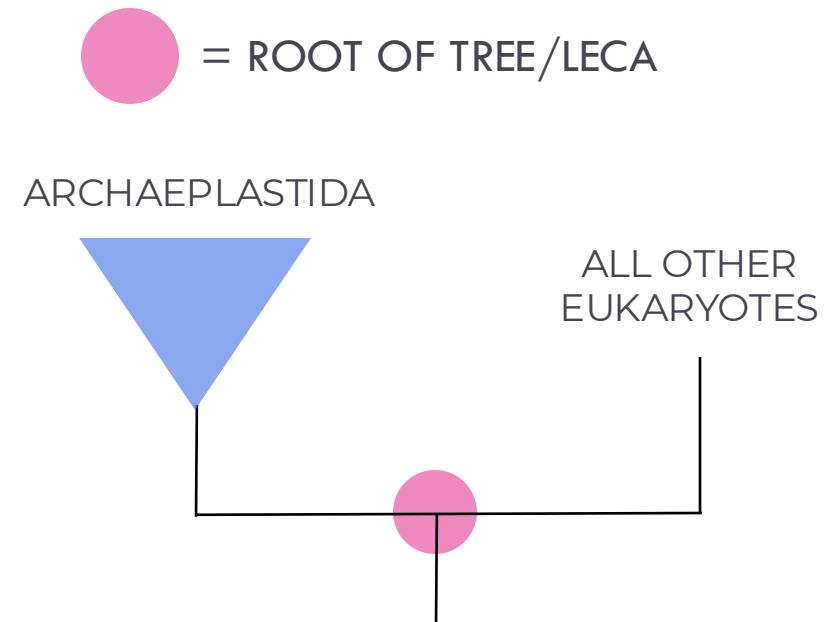


TESTING ALTERNATE ROOT POSITIONS

Alternative root topologies were generated by setting monophyly constraints in IQ-TREE

Additional roots tested:

1. **METAMONADA (ANAERAMOEBAE)**
2. **DISCOBA**
3. **OPISTHOKONTA**
4. **MALAWIMONADA**
5. **JAKOBIDA**
6. **EUGLENOZOA**
7. **ARCHAEOPLASTIDA** ←



INCLUDING ADDITIONAL COMPLEXITY TO THE PHYLOGENETIC MODELS

- FUNDI** – models **functional divergence** (change of amino acid preference) at sites across a known branch (ie, the branch between alpha and euka)
- HOST** – models **heterotachy** (changing rate of evolution at sites across the tree)
- GF-MIX** – accommodates **amino acid compositional biases** driven by changing GC content across the tree

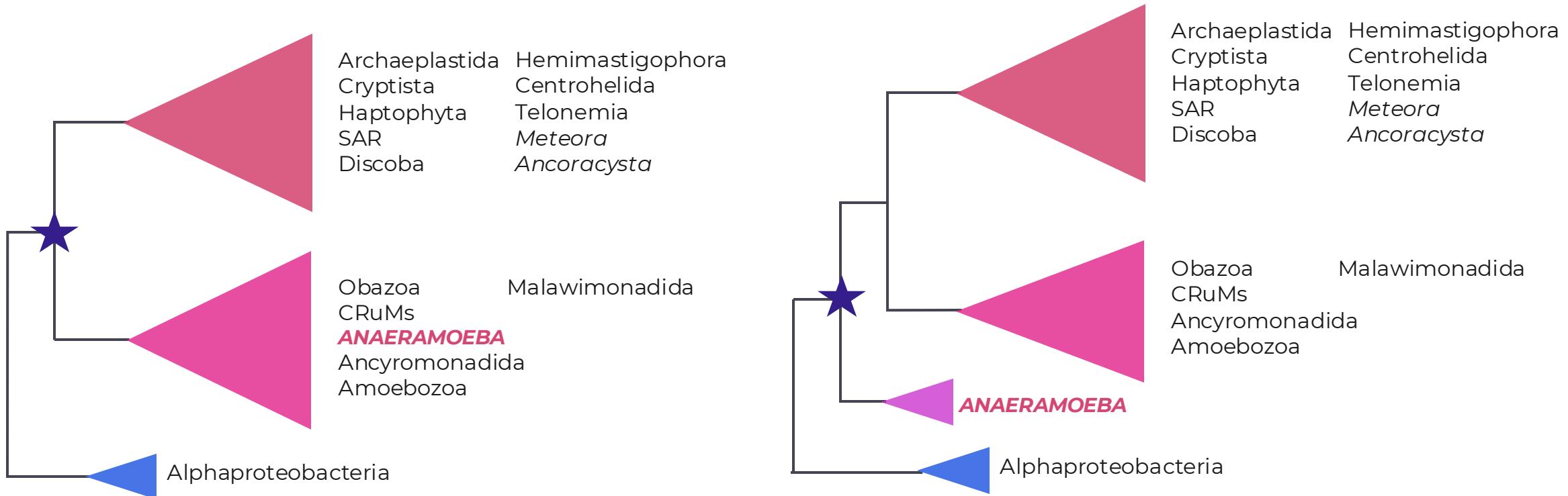
Likelihoods **improve** under more complex models, but order of root preference remains consistent.

Neither an Opimoda+ root or an Anaeramoeba root are rejected by topology testing.

ANAE+ DATASET				
TOPOLOGY	MEOW80	MEOW80+FUNDI	MEOW80+GHOST	MEOW80+GF-MIX
OPIMODA+	-1665043.854	-1657618.898	-1650367.99	-1655099.392
ANAERAMOEBA	-1665044.498	-1657619.576	-1650369.352	-1655103.092
DISCOBA	-1665061.809	-1657635.485	-1650394.794	-1655134.492
MALAWIMONADIDA	-1665095.632	-1657650.686	-1650414.421	-1655148.935
OPISTHOKONTA	-1665138.368	-1657687.757	-1650454.564	-1655204.365
EUGLENOZOA	-1665146.579	-1657718.133	-1650475.147	-1655210.287
ARCHAEPLASTIDA	-1665170.708	-1657713.761	-1650467.552	-1655221.351
JAKOBIDA	-1665205.285	-1657772.885	-1650538.914	-1655281.622

The position of *Anaeramoeba*, the only representative of Metamonada, is not well-resolved

- ANAEROBES – NO MITOCHONDRIAL GENOME
- SITE OCCUPANCY OF 18%
- LONG-BRANCHING

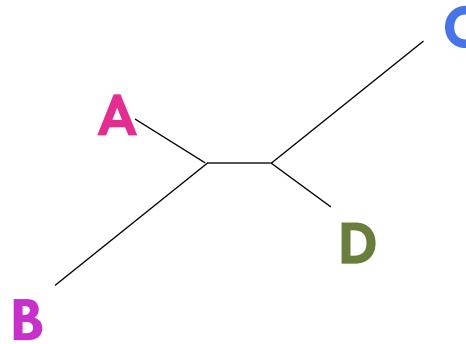


★ = ROOT

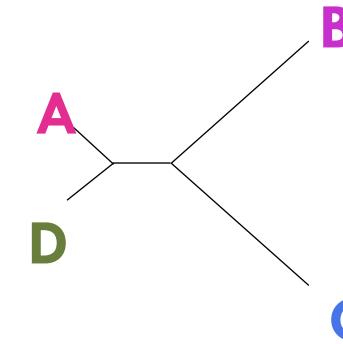
LONG BRANCH ATTRACTION (LBA)

Under model misspecification or small sample bias, long branches artefactually branch together

TRUE TREE

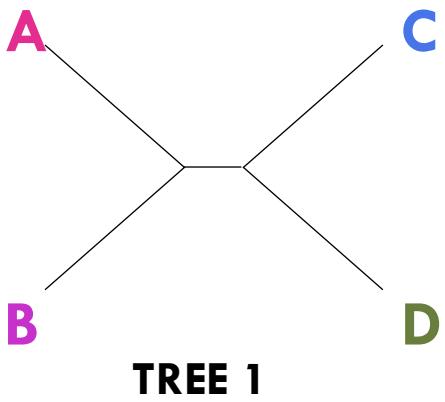
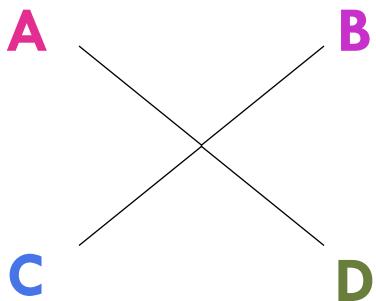


LBA TREE

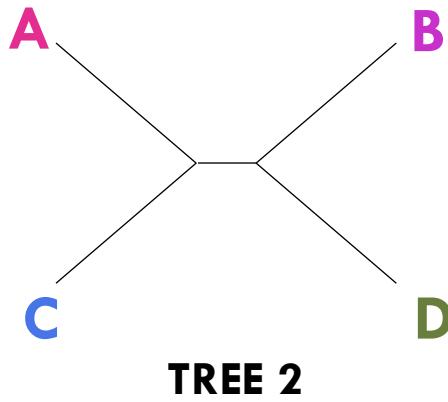


Simulation studies support an artefactual attraction of *Anaeramoeba* to the long branch connecting the outgroup

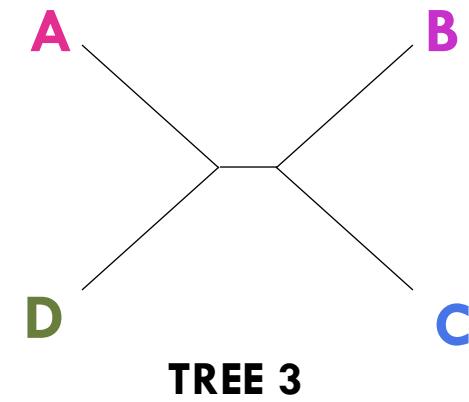
TRUE TREE



TREE 1



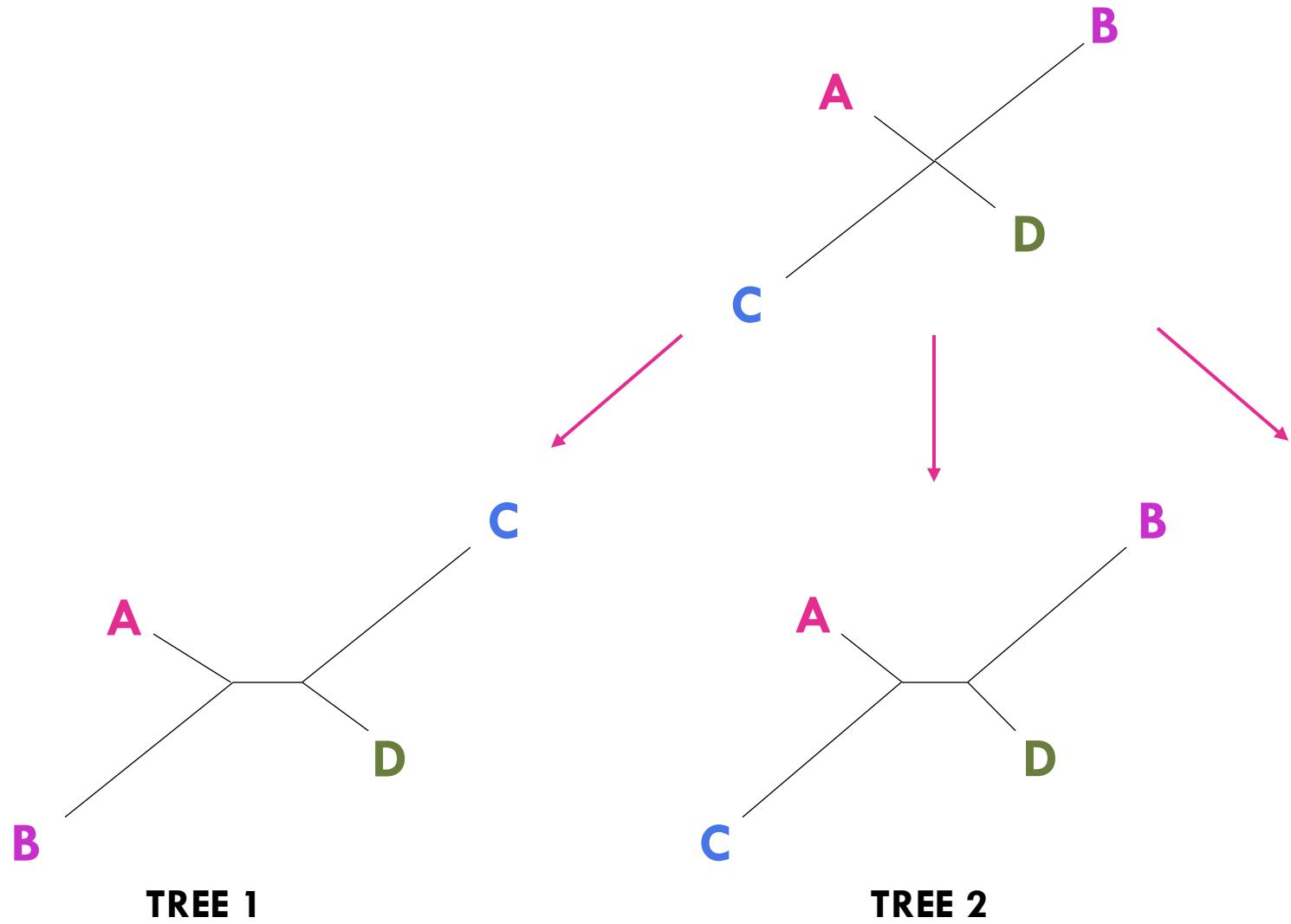
TREE 2



TREE 3

ALL TOPOLOGIES EQUIALLY LIKELY

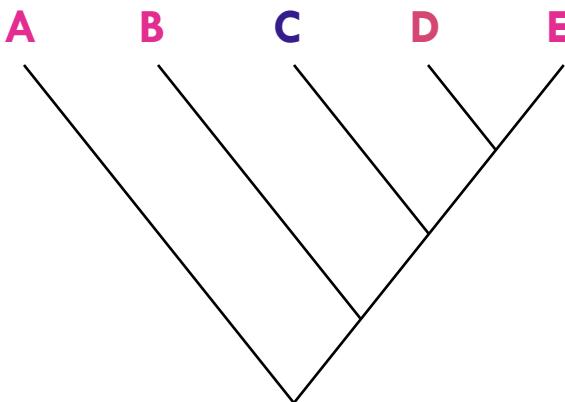
TRUE TREE



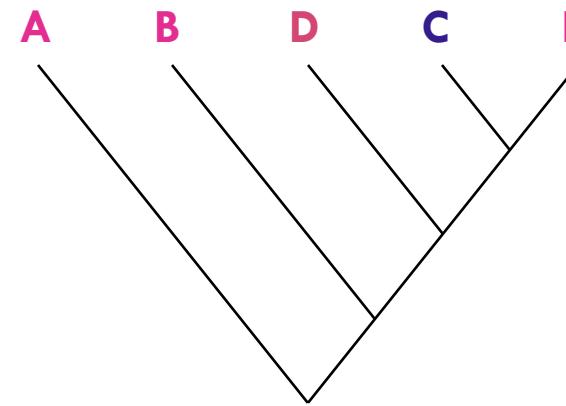
UNDER MODEL MISSPECIFICATION OR SMALL SAMPLE BIAS, LBA TREE IS PREFERRED

TESTING FOR LONG BRANCH ATTRACTION

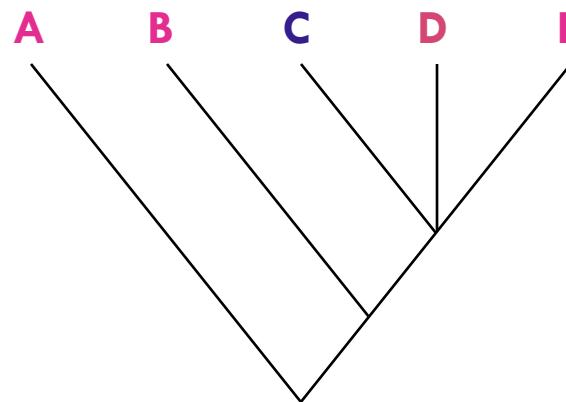
1. Create a consensus tree from two competing topologies



VS.



CONSENSUS:



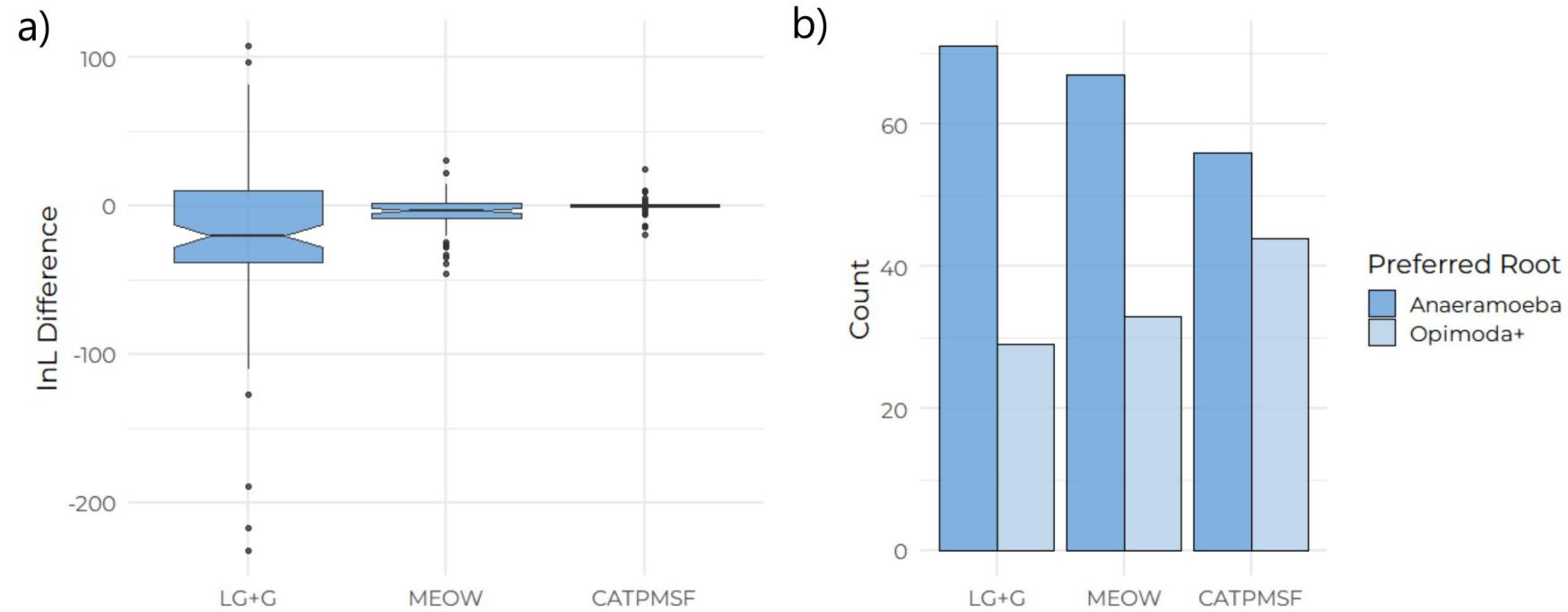
TESTING FOR LONG BRANCH ATTRACTION

1. Create a consensus tree from two competing topologies – **consider this the ‘true’ topology**
2. Simulate 100 alignments based on the consensus tree and CAT-PMSF site profiles (ALISIM)
3. Run IQ-TREE on all alignments to calculate the likelihood for both resolved topologies (Opimoda-like and *Anaeramoeba*) given the following models:
 1. CAT-PMSF
 2. MEOW
 3. LG+G

If there is **no LBA**, there should be no preference for one topology over the other (ie. Both topologies are equally likely to have the better likelihood in a given replicate)

If there is **LBA**, there will be a bias toward the *Anaeramoeba* root topology under model misspecification (ie. *Anaeramoeba* topology will have **better** likelihood)

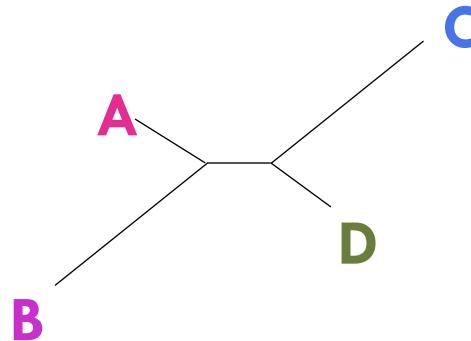
As model misspecification increases, preference for an *Anaeramoeba* root over an *Opimoda+* root also increases



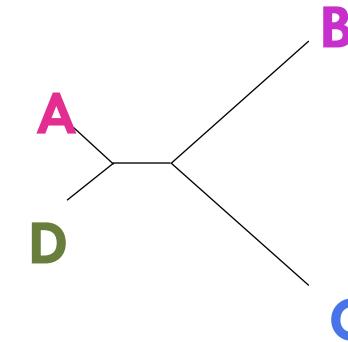
LONG BRANCH ATTRACTION (LBA)

Under model misspecification or small sample bias, long branches artefactually branch together

TRUE TREE



LBA TREE



Simulation studies support an artefactual attraction of *Anaeramoeba* to the long branch connecting the outgroup

→ *Anaeramoeba* is sensitive to LBA – not enough information to place them confidently in the tree

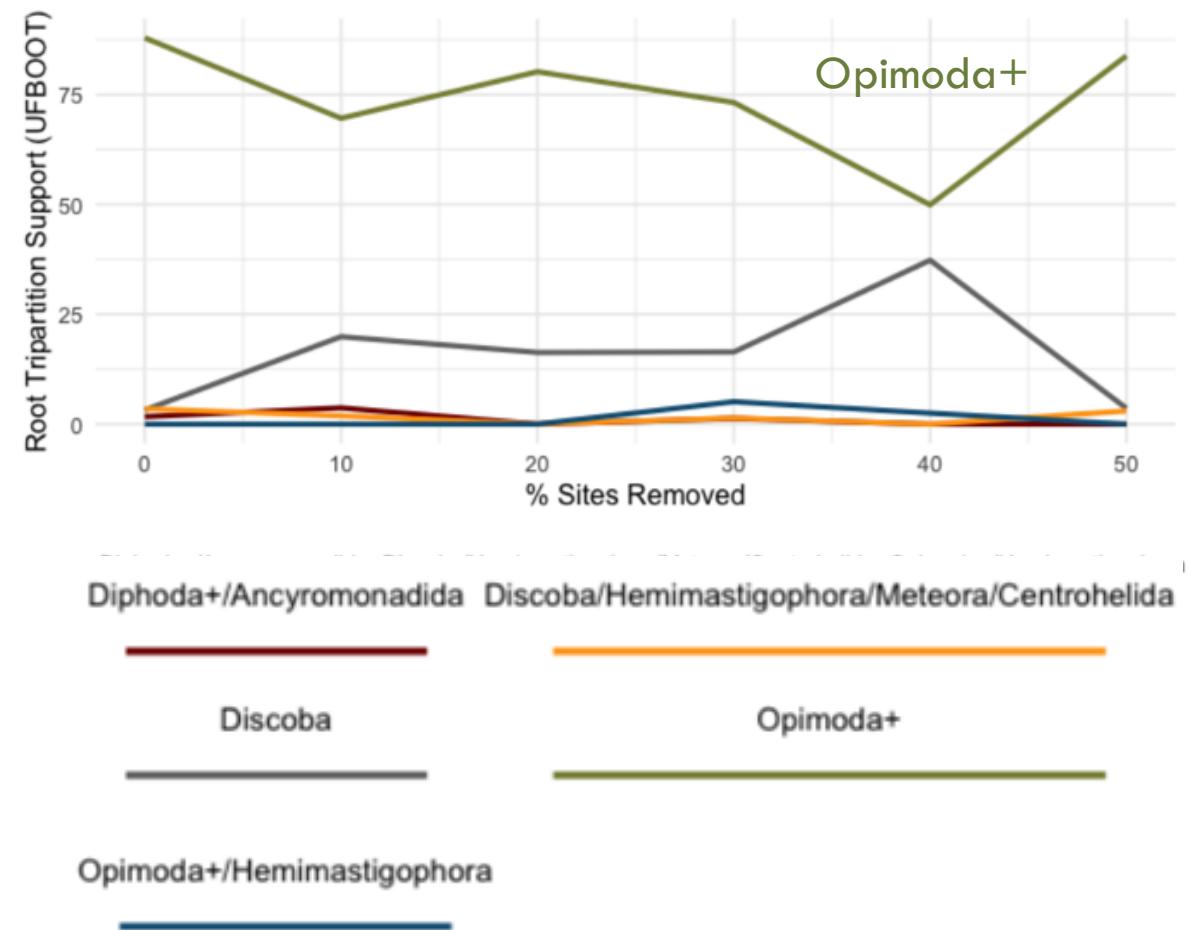
Removing the fastest-evolving sites, genes, and taxa

REMOVAL OF FASTEST EVOLVING SITES

A

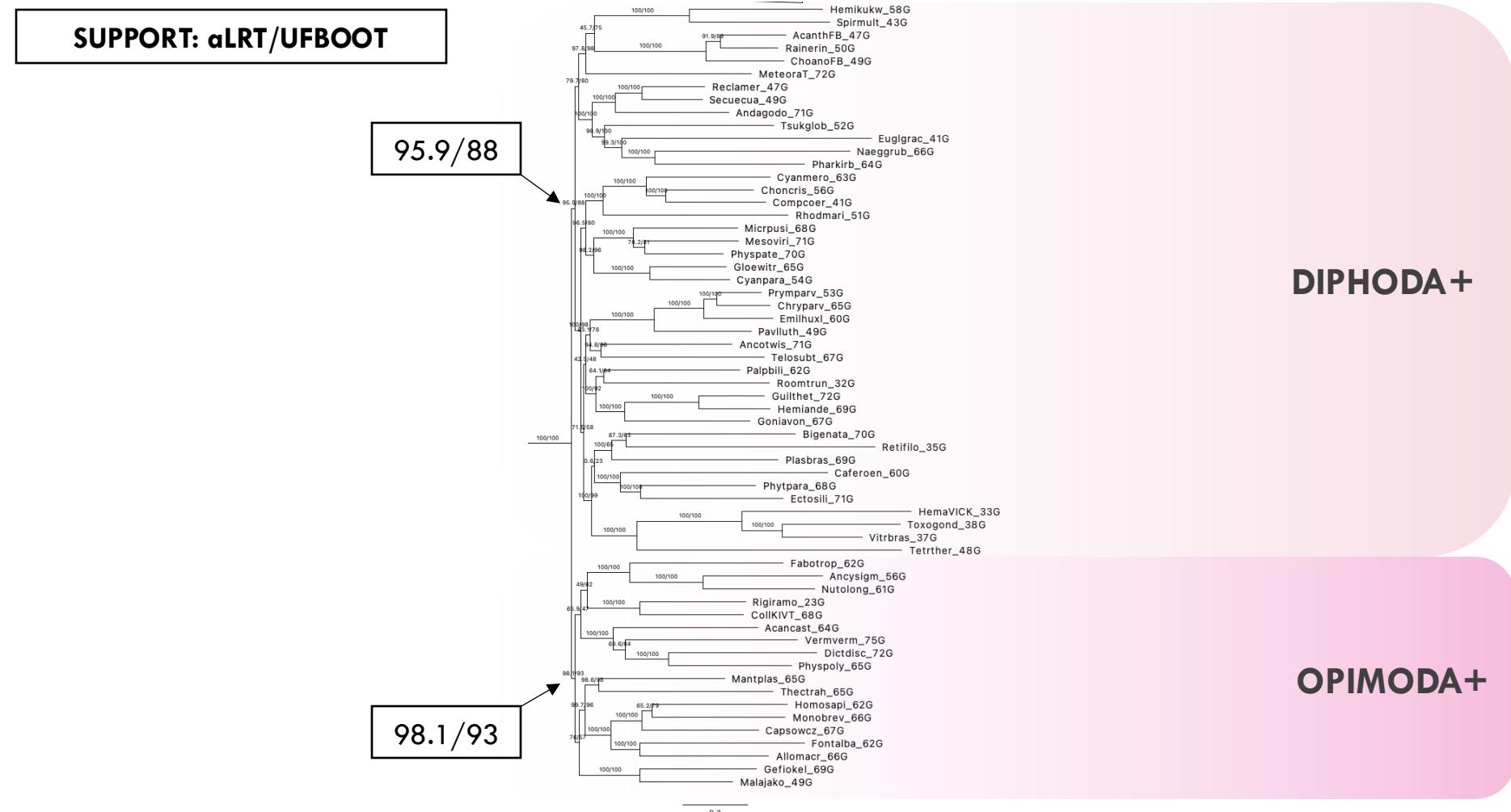
% Sites Removed	aLRT/UFBOOT Support	
	Opimoda+	Diphoda+
10	98.5/93	53.2/75
20	99.4/98	78.7/82
30	99.2/93	50.3/75
40	94.6/93	30.2/55
50	98.0/97	95.8/87

B



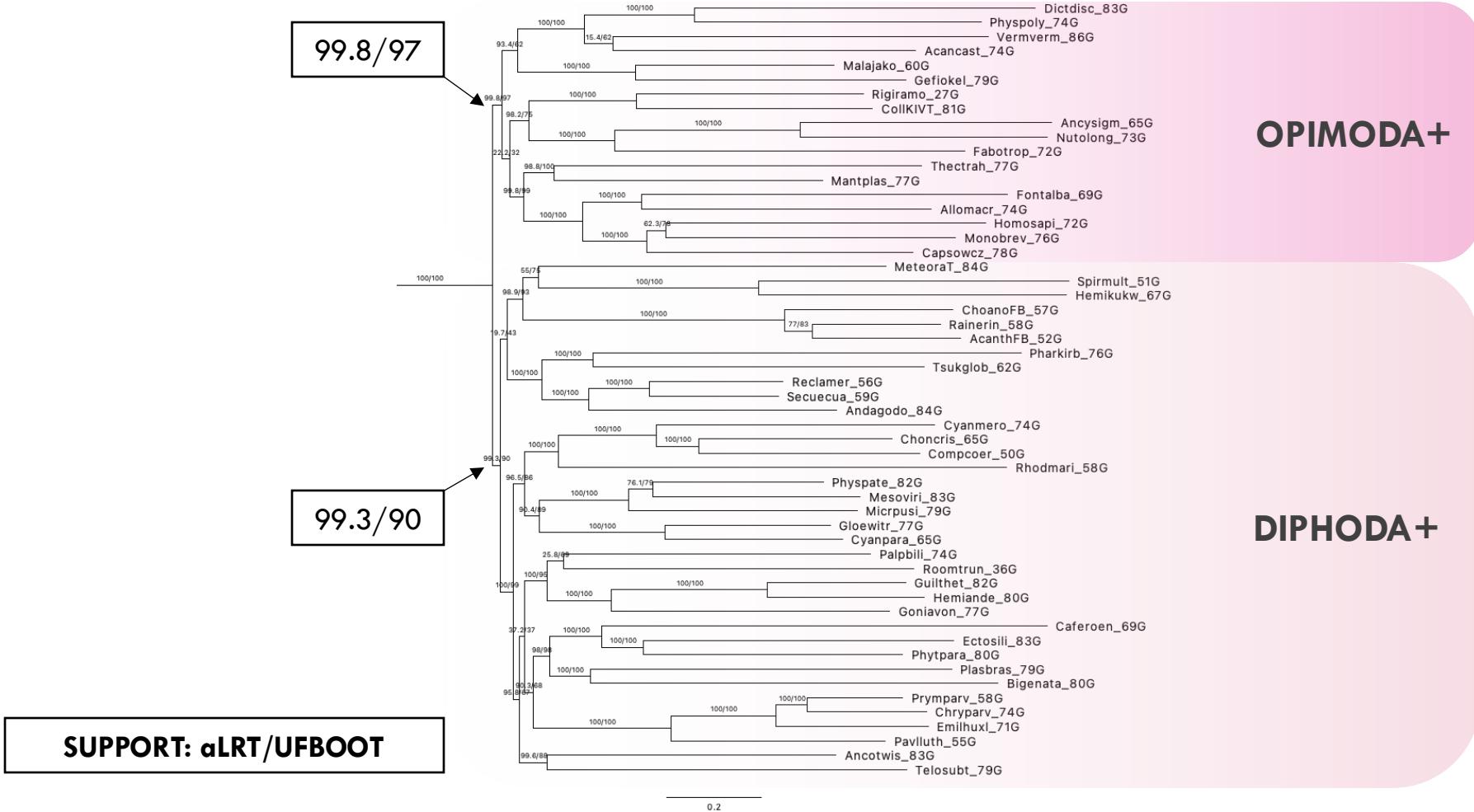
REMOVAL OF 13 MOST DIVERGENT GENES

Based on the internal branch length connecting eukaryotes and Alphaproteobacteria

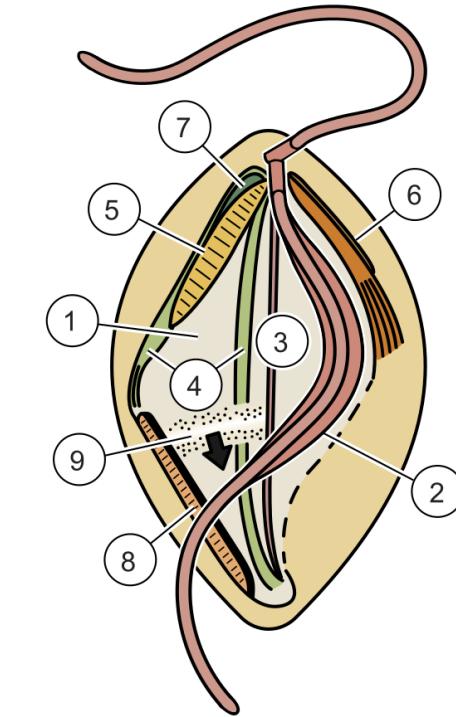
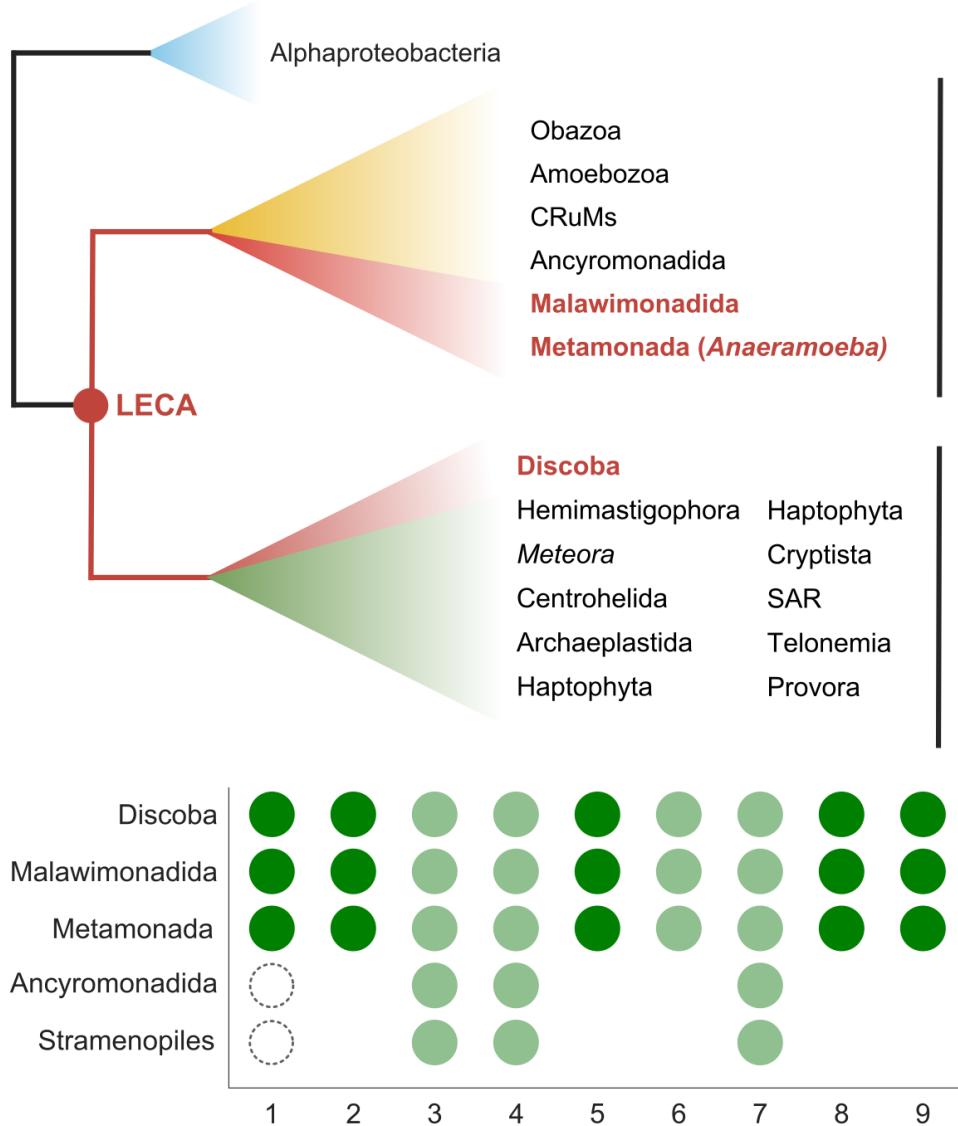


REMOVAL OF 7 FASTEST EVOLVING TAXA

Based on tip-to-tip distance



The recovered root position suggests LECA was an excavate-like organism

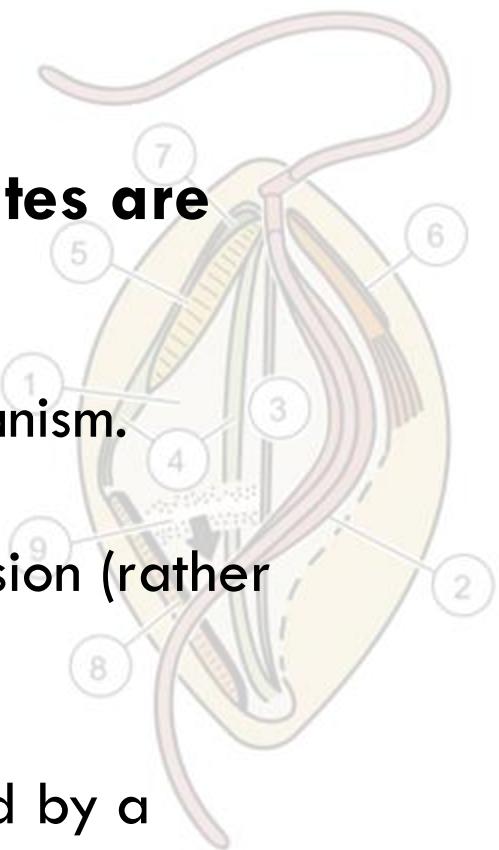


SUMMARY

- **New, much larger mitochondrial protein dataset** for rooting the eukaryote tree
- **All tested models recover an ‘Opimoda+’ root** in both maximum likelihood and Bayesian frameworks
- Root position is **robust to the removal of the fastest evolving sites, genes, and taxa**

Complex cytoskeletal elements common to ‘typical’ excavates are features of LECA that were lost in other lineages

- LECA was likely a small (25 µm or less), unicellular, and flagellated organism.
- Phagotrophic (likely bacterivorous), probably feeding on prey in suspension (rather than as a surface feeder like many amoebae).
- Had a defined cell shape crucial to its motility and feeding, underpinned by a complex cytoskeletal organisation, including microtubular roots of diverse sizes, functions, and associated non-microtubular elements.
- **Models of eukaryogenesis need to have such a form of eukaryote as their “endpoint”.**



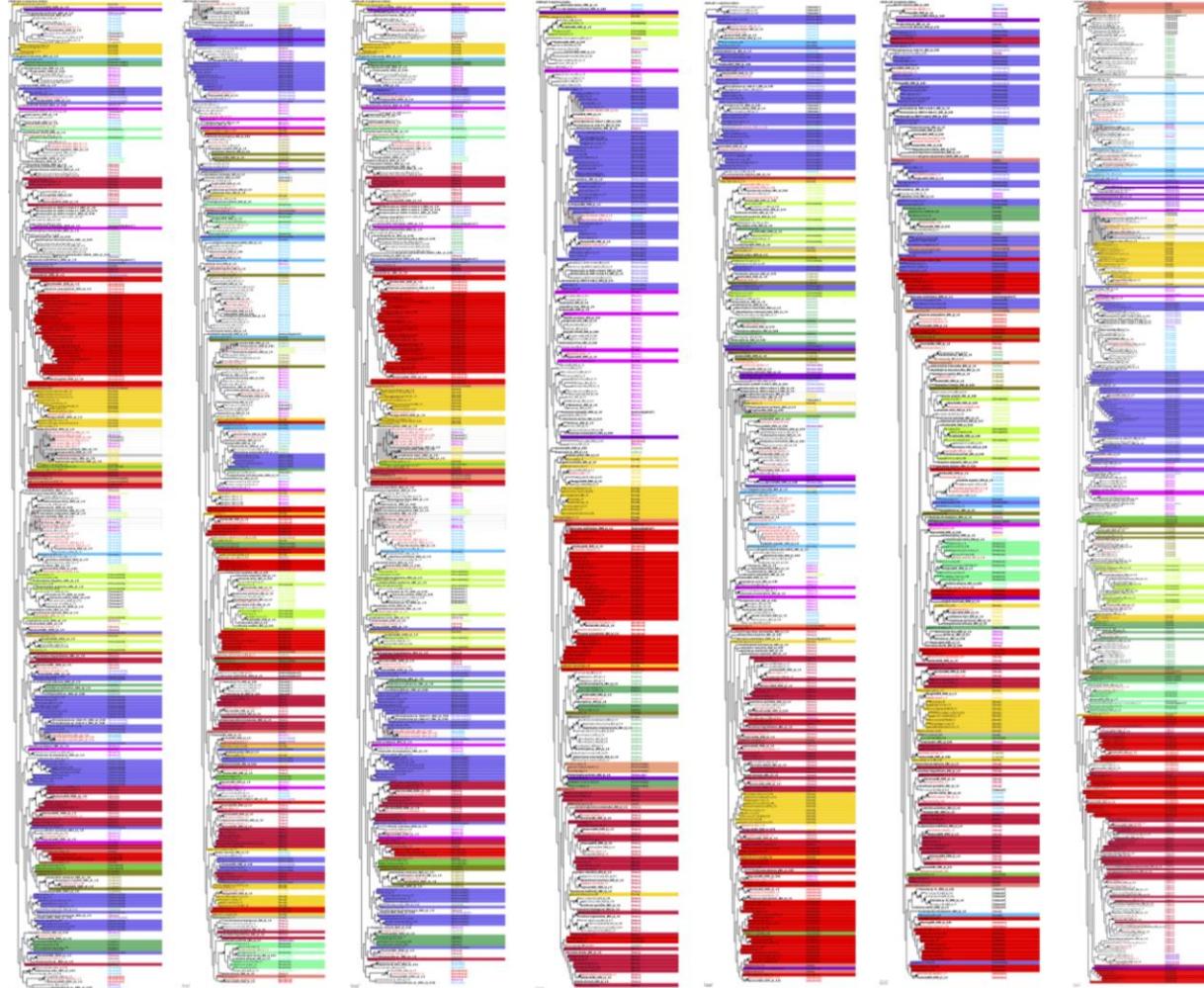
Useful tools

PHYLOFISHER

A phylogenetically aware pipeline for phylogenomic dataset construction

David Žihala, Alexander K. Tice, Tomáš Pánek, Serafim Nenarokov, Eric Salomaki,
Andrew J. Roger, Martin Kolísko, Fabien Burki, Laura Eme, Marek Eliáš,
Matthew W. Brown

Selection of orthologs and orthologous sequences in them is critical

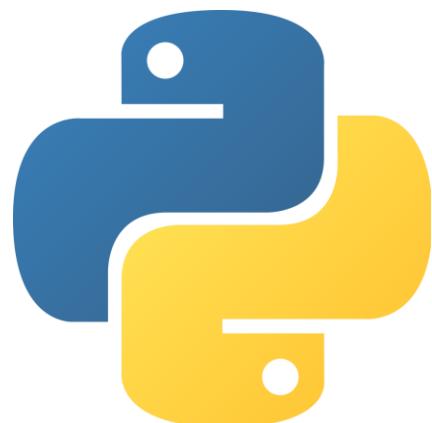


- Contamination
 - On sequencer
 - Endosymbionts
 - Prey (or predators)
- Paralogs
 - Genomic duplication
 - Deep- (i.e., a- vs b-tubulin)
 - Mid- (within a group)
 - In- (within a species (or genus))
- Phylogenetically informative
 - Broad taxonomic sampling
- To do this requires trees and careful consideration of them
 - Eyes

PHYLOFISHER

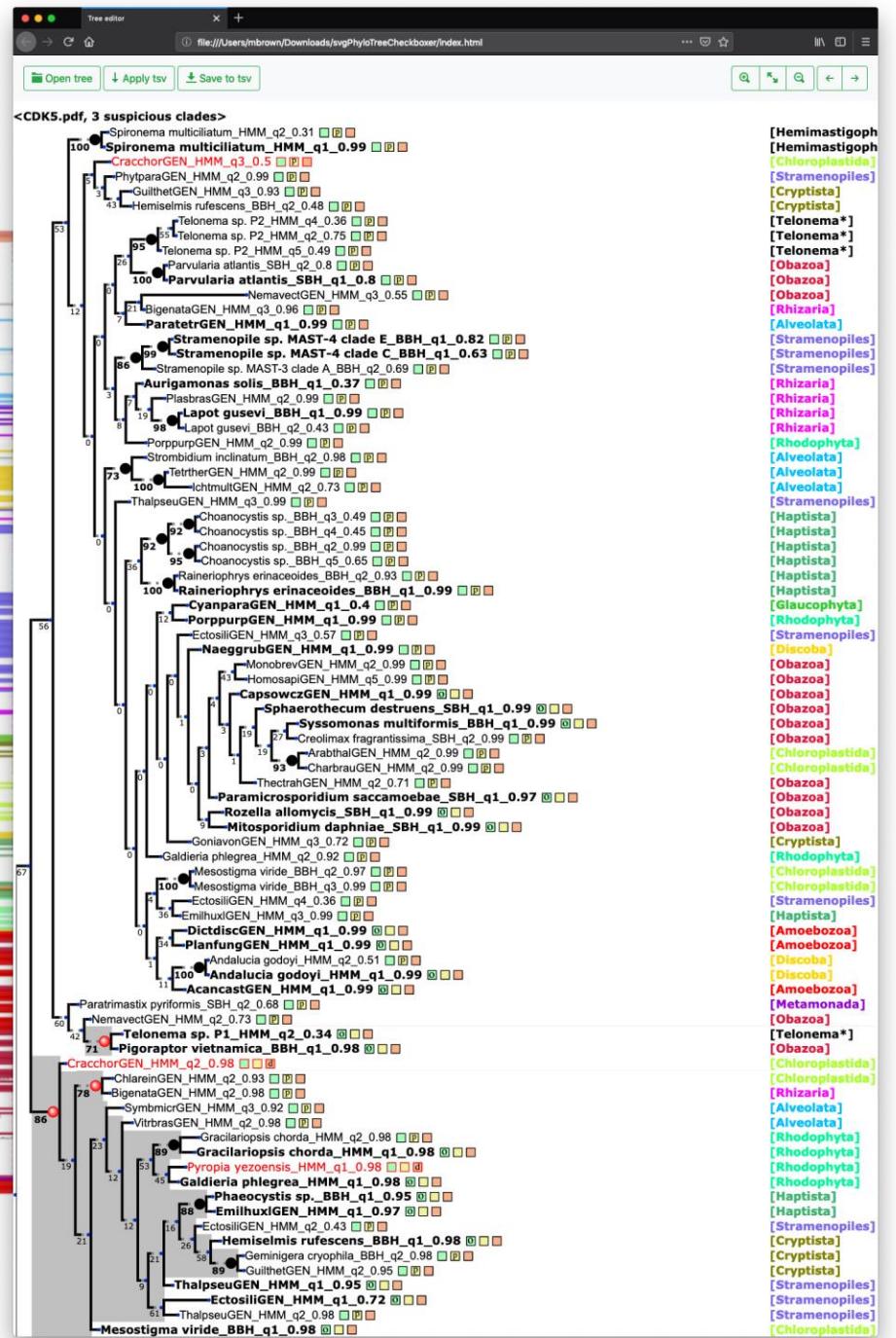
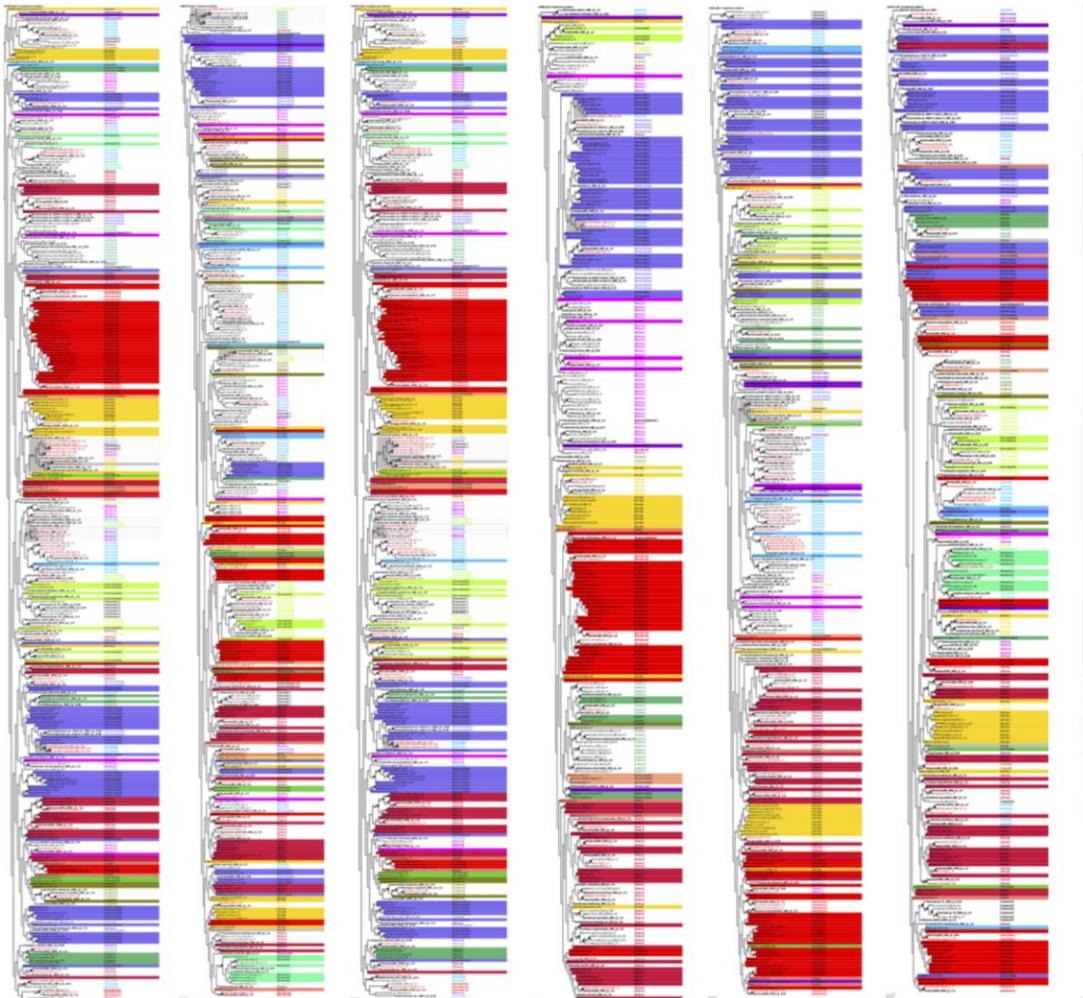
New method (and tool) allows for others
to simply do phylogenomics

- Ships with a phylogenomic matrix and tool (via GitHub)
 - 310 Taxa, covering all deep eukaryotic groups
 - 240 Orthologs (whittled down from Brown et al. 2018)
- Coded in Python in a easy to to install CONDA environment
 - All dependencies are automatically installed

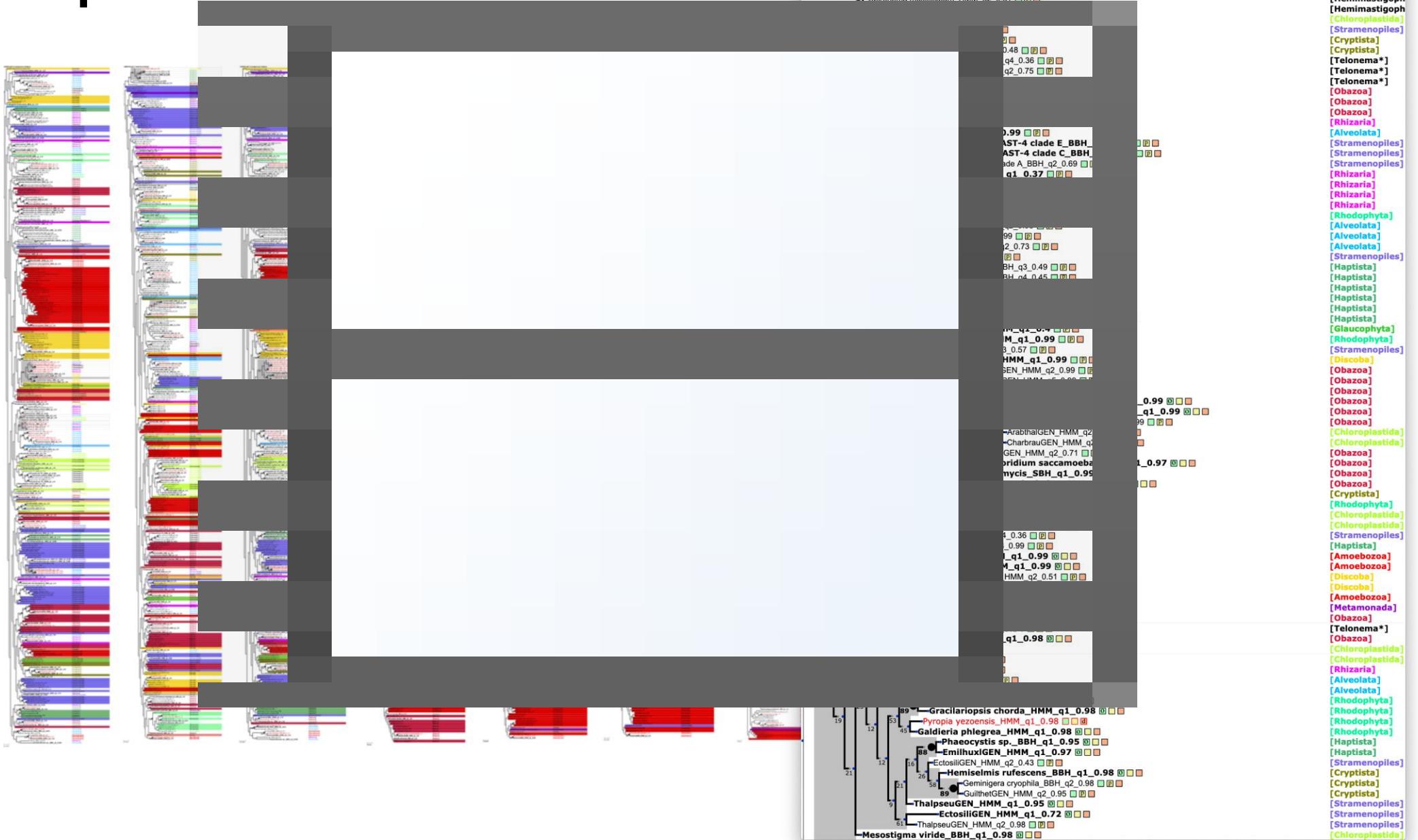


conda env create -f fisher_env.yml

Tree Inspection

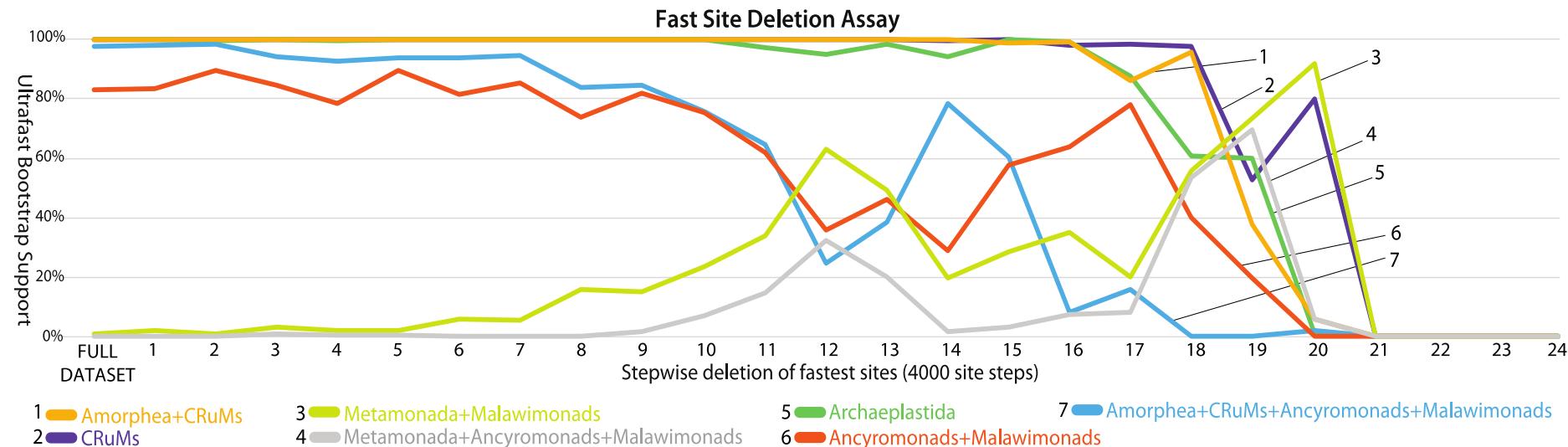


Tree Inspection



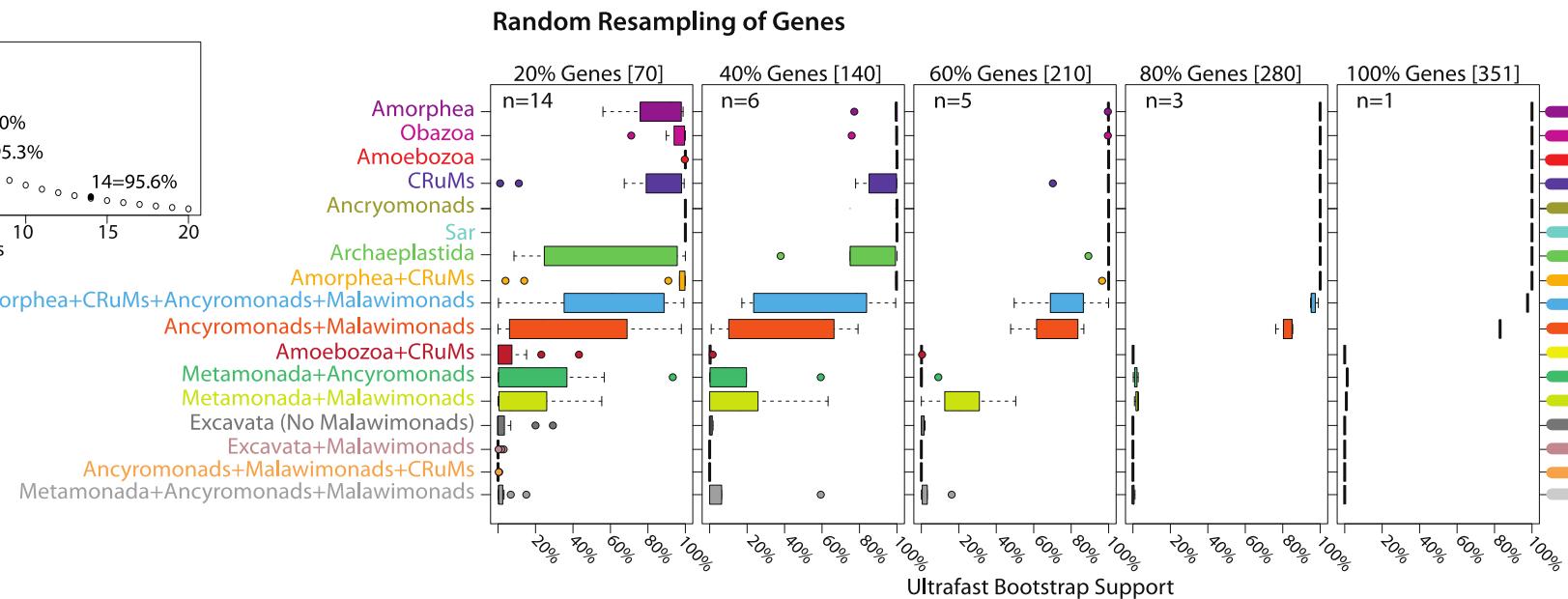
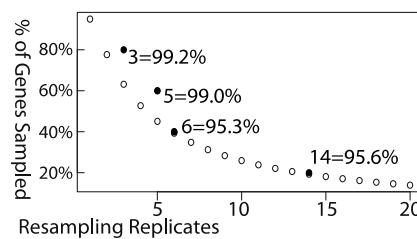
PHYLOFISHER

- Easily installed and simple usage
- Ships with our dataset
 - Includes Paralogs for tree building, more accurate identification
- Your own gene sets can be incorporated or used independently
- Tools for post-phylogenomic analyses

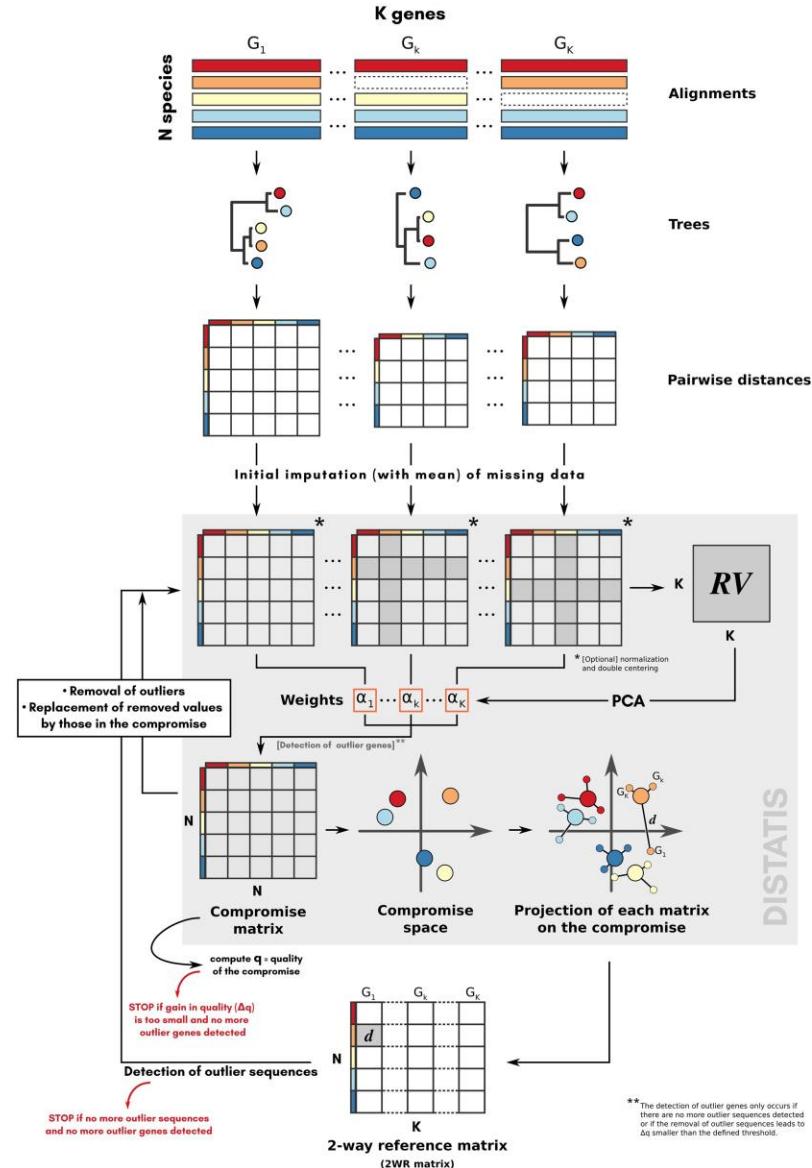


PHYLOFISHER

- Easily installed and simple usage
- Ships with our dataset
 - Includes Paralogs for tree building, more accurate identification
- Your own gene sets can be incorporated or used independently
- Tools for post-phylogenomic analyses



Phylter

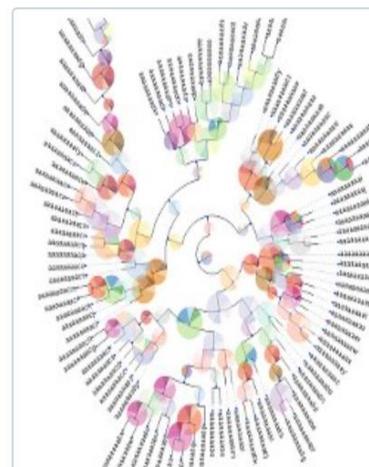


ETE3

A Python framework for the analysis and visualization of trees.

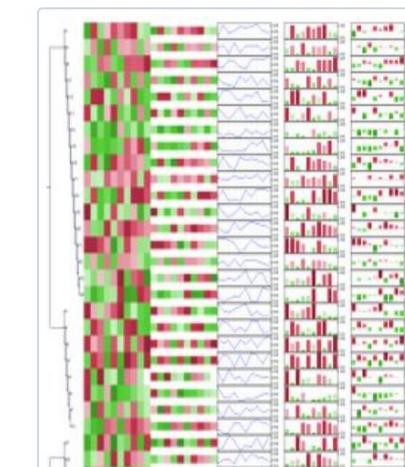
[Download](#)[Python API](#)[Cookbook](#)[Phylogenomic tools](#)[Contribute](#)

```
from ete3 import Tree
tree = Tree('((A,B), D);')
print tree
#      /-A
#      /-|
#      -|- \-B
#      \-D
A = tree & "A"
A.up.show()
```



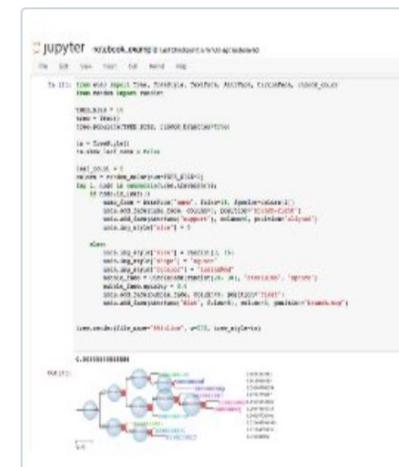
Trees as Python objects

Load, create, traverse, search, prune, or modify hierarchical tree structures with ease using the ETE Python API.



Programmatic tree visualization

Get full control of your tree images. Browse them interactively or render SVG, PNG or PDF images.



Tree annotation

Custom node attributes can be rendered as graphical elements. Choose among external images, charts, symbols, text labels, and

Jupyter notebook support

Prototype your methods using the Jupyter notebook framework including inline visualization of trees.