

Estimation of evolutionary distance based on K -string composition

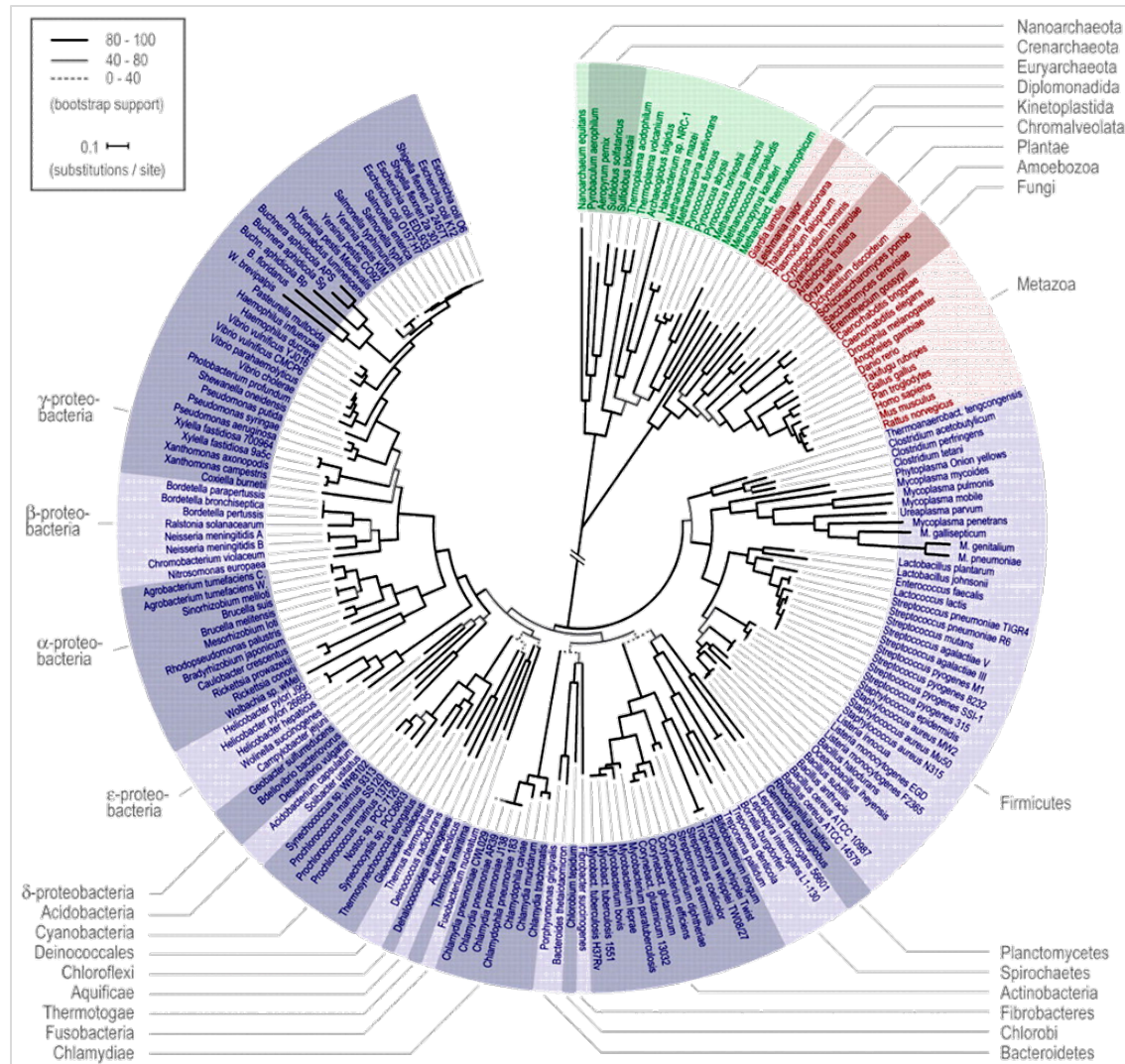
Qiang Li

Department of Combinatorics and Geometry
CAS-MPG Partner Institute for Computational Biology



Difficulties in phylogenomics

- Species trees versus gene trees
- Too few common genes
 - “The tree of one percent”
- Model selection
 - A science, or an art?
- Computational complexity
- ...



A tree based on 31 genes (Ciccarelli '06)

Phylogenetics above sequence-level

- Rare genomic change
- Gene content & gene order
- *K*-string composition

“In summary, the methods for inferring trees based on whole-genome features are at an early stage of their development, which might be comparable to that of sequence-based methods in the early 1970s. In particular, they generally lack a global probabilistic modeling.”

(Philippe 2005)

CVTree: Compositional vector

- “Random background”: $(K - 2)$ -order Markov chain

$$f^0(\alpha_1 \cdots \alpha_K) = \frac{f(\alpha_1 \cdots \alpha_{K-1}) f(\alpha_2 \cdots \alpha_K)}{f(\alpha_2 \cdots \alpha_{K-1})}$$

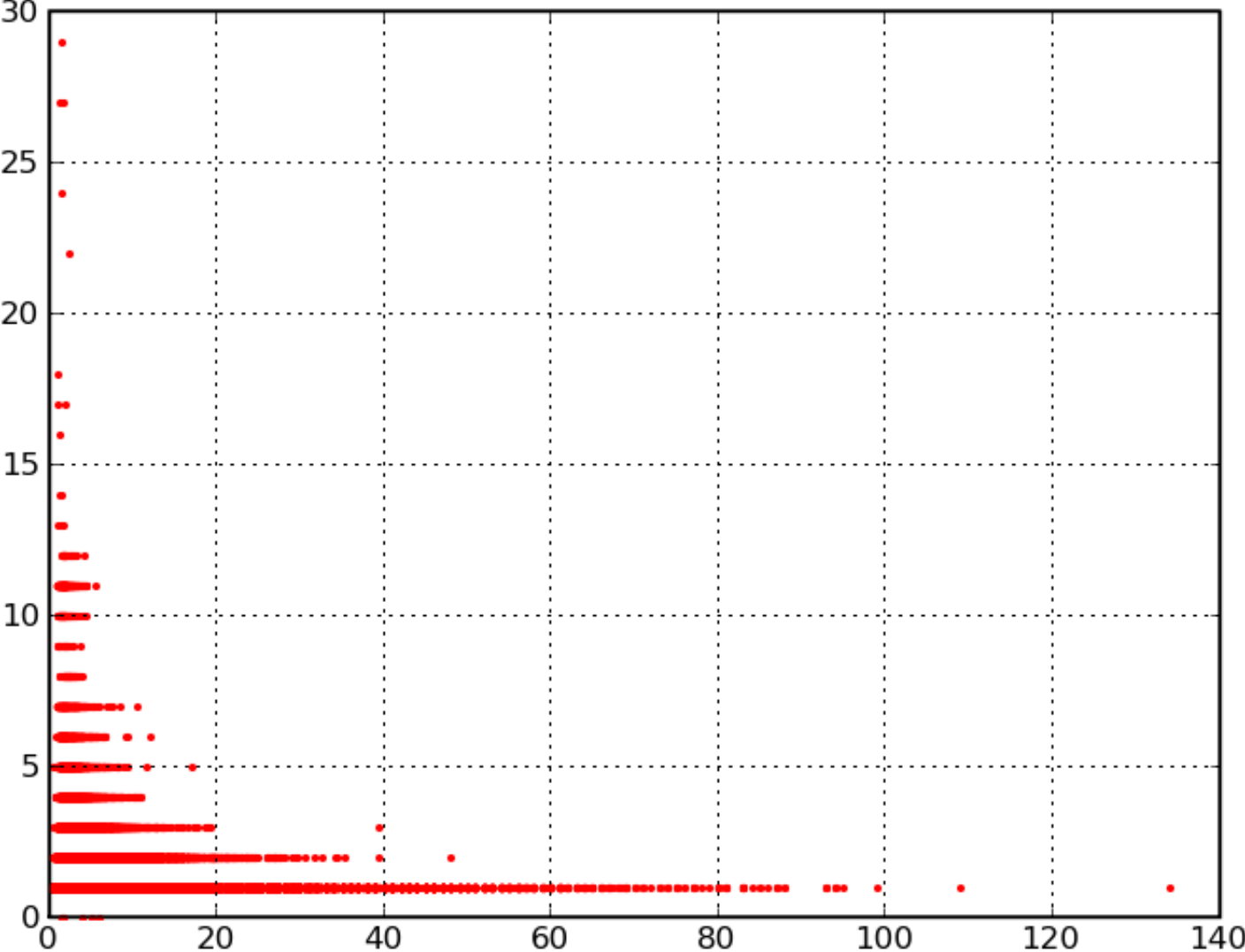
- “Subtraction”

$$a_i = \begin{cases} \frac{f(i) - f^0(i)}{f^0(i)}, & f^0(i) \neq 0 \\ 0, & f^0(i) = 0 \end{cases} \quad i \in \Sigma^K$$

- Distance (dissimilarity)

$$d(\mathbf{a}, \mathbf{b}) = \frac{1 - \cos(\mathbf{a}, \mathbf{b})}{2} \in [0, 1]$$

eco6_CV_occur_plottable

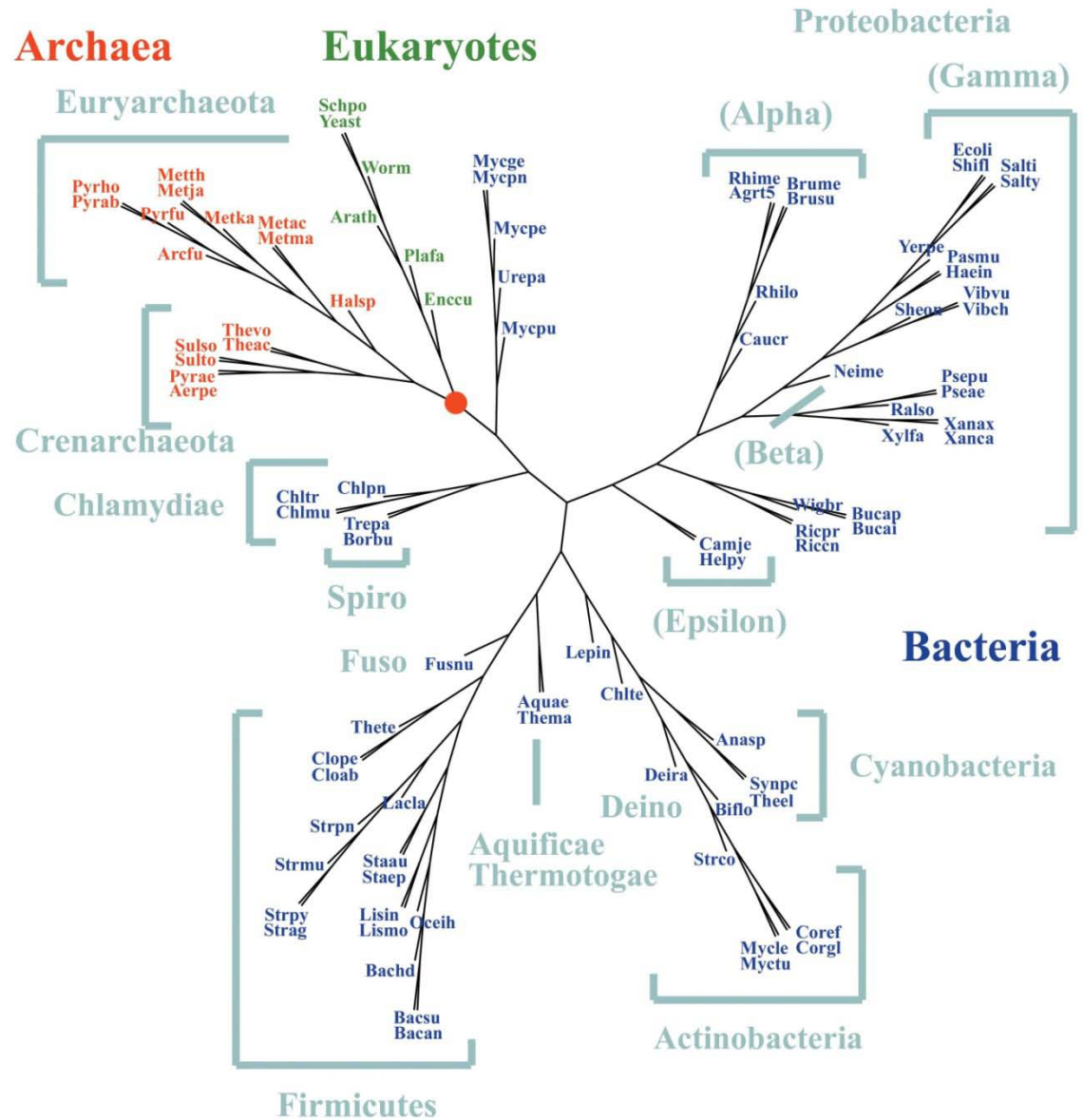


f versus f^0

Courtesy of Alexander Grossmann

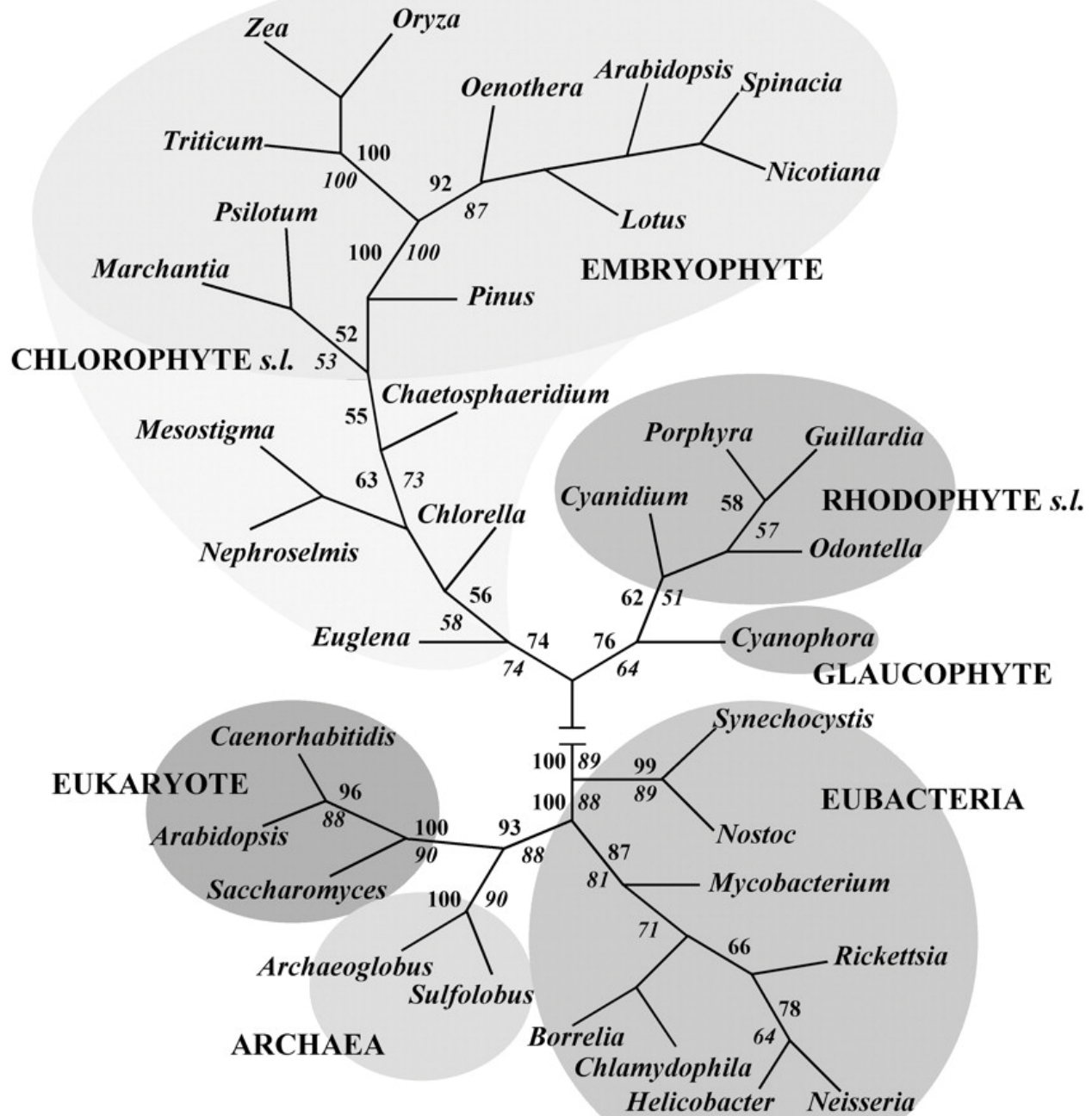
CVTree

- simple
- efficient
- parameter free
- truly whole genomic
- ...
- mysterious

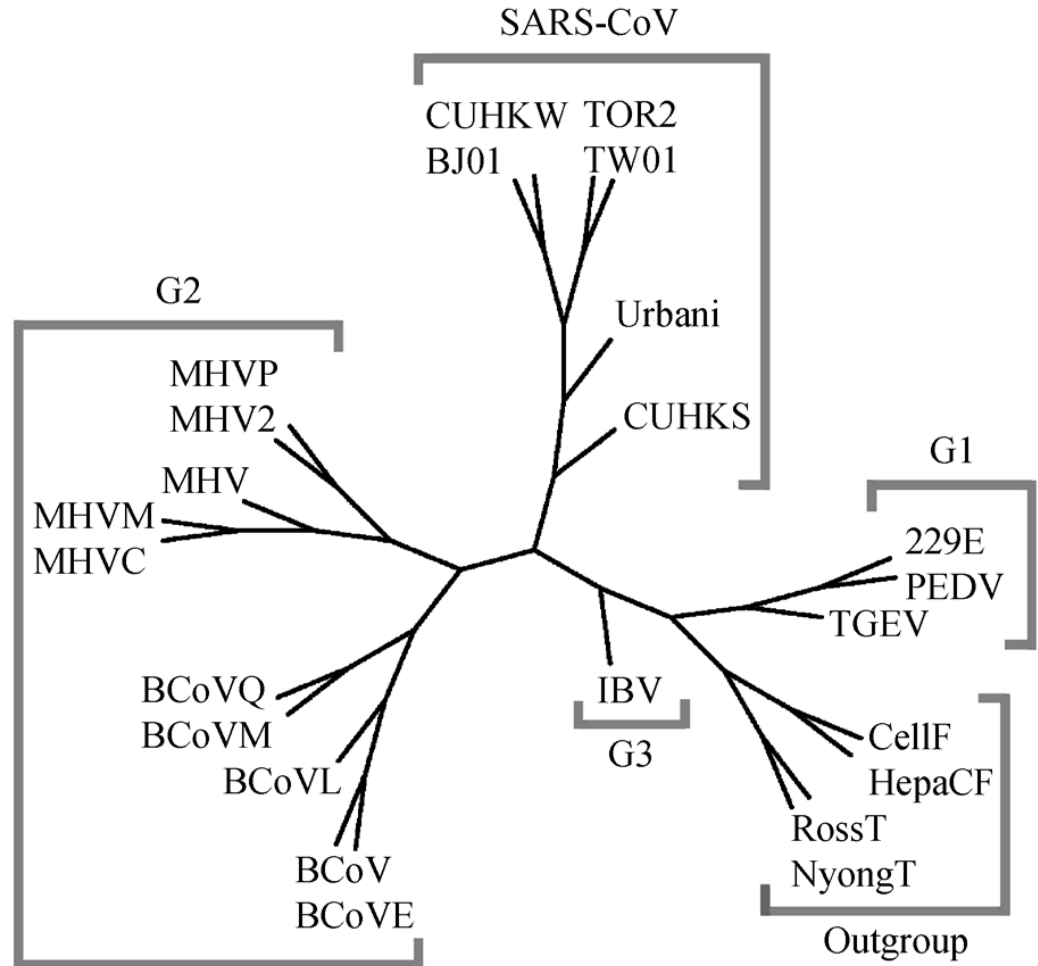


CVTree of chloroplasts

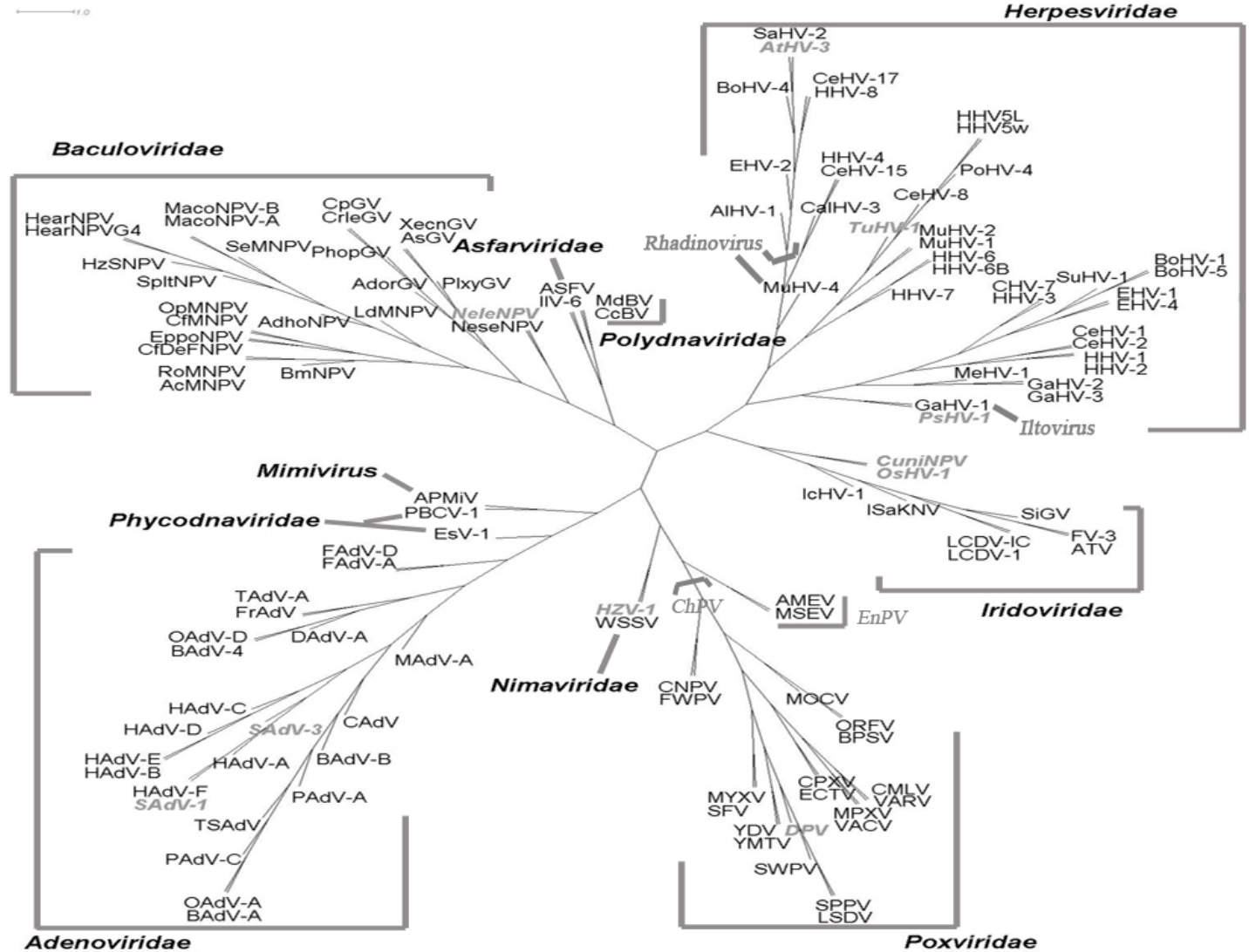
a.



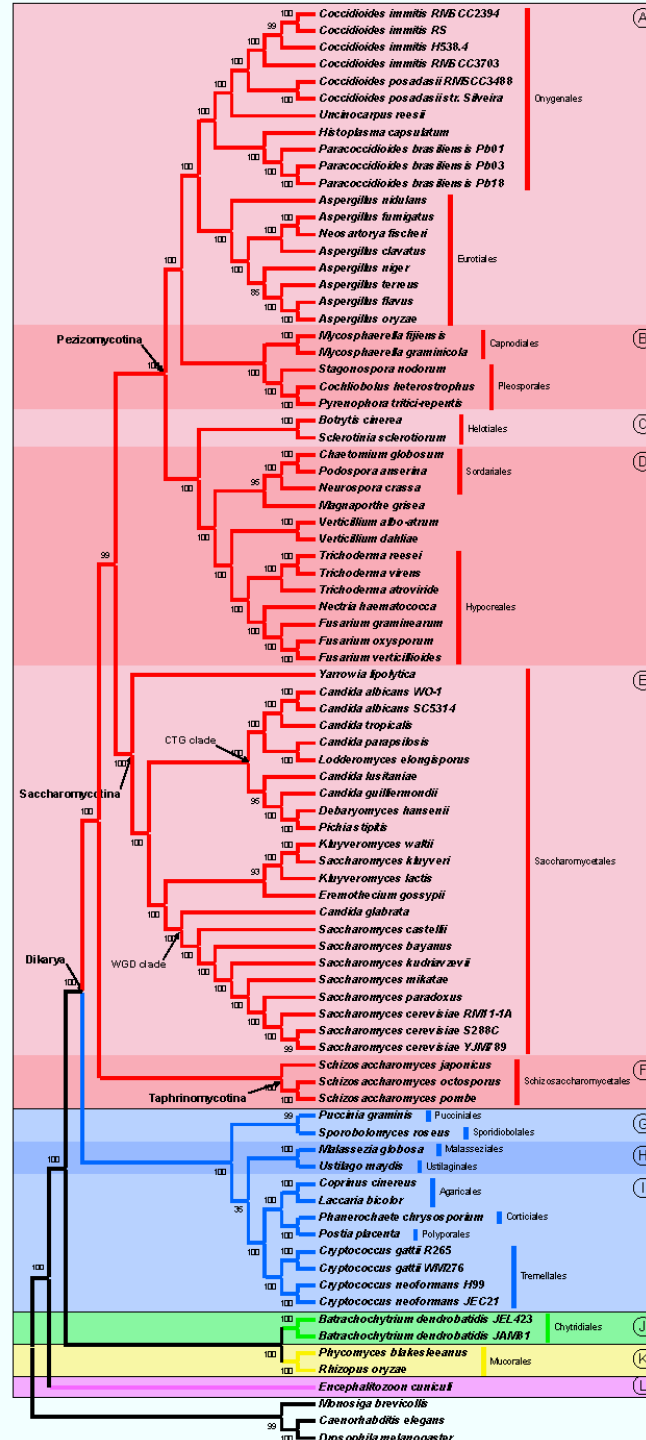
CVTree of coronavirus

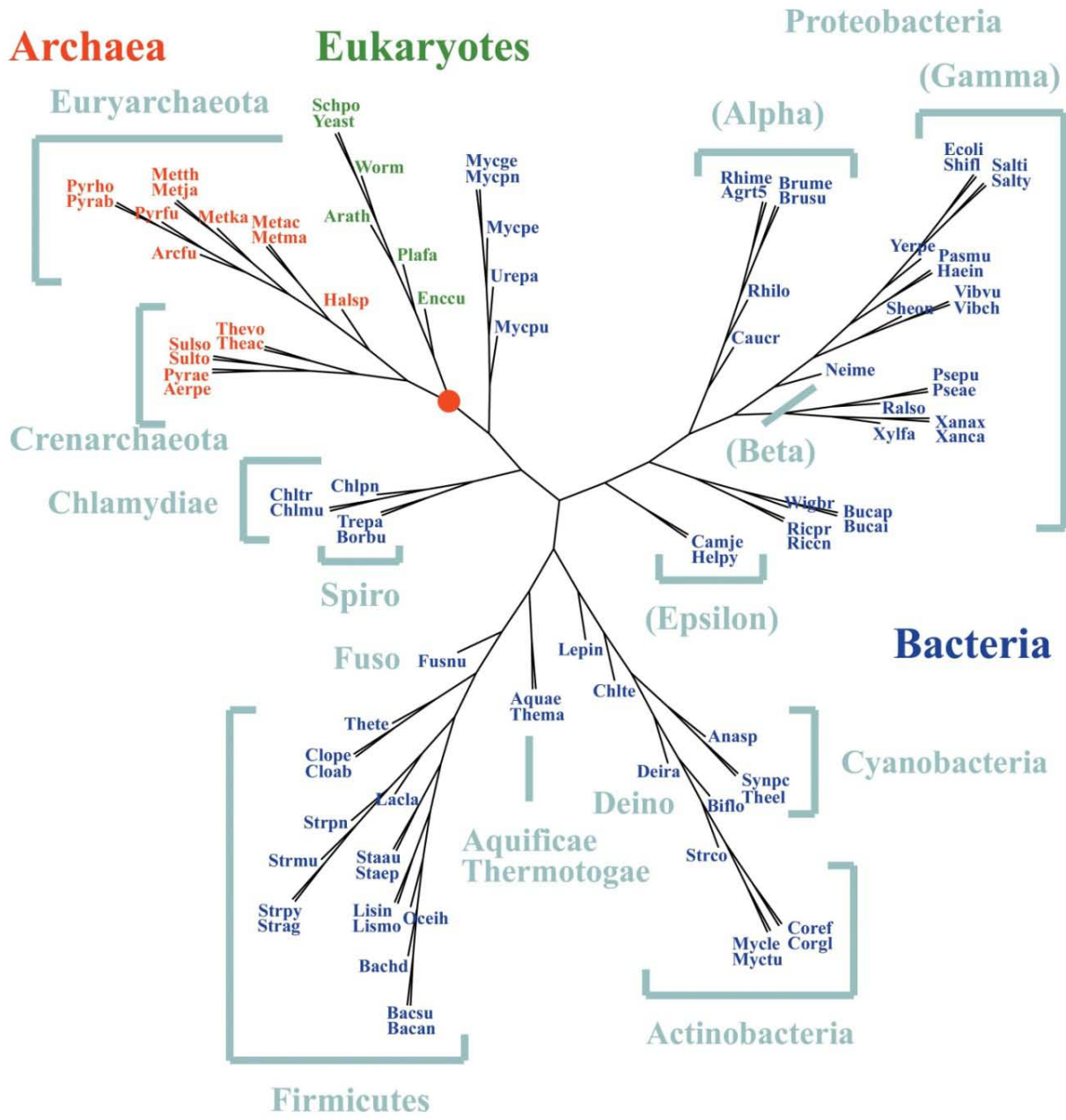


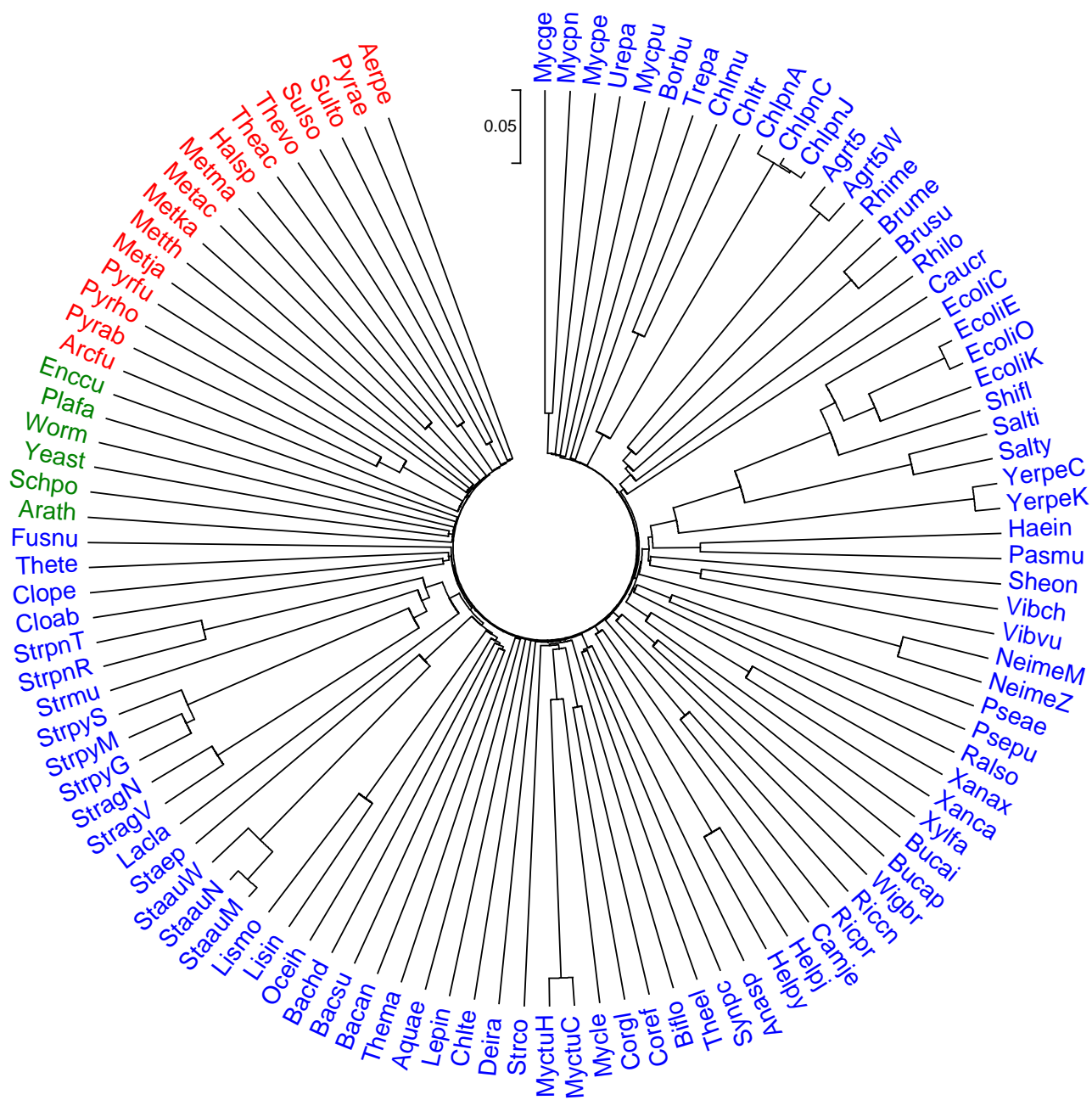
CVTree of dsDNA viruses



CVTree of fungi



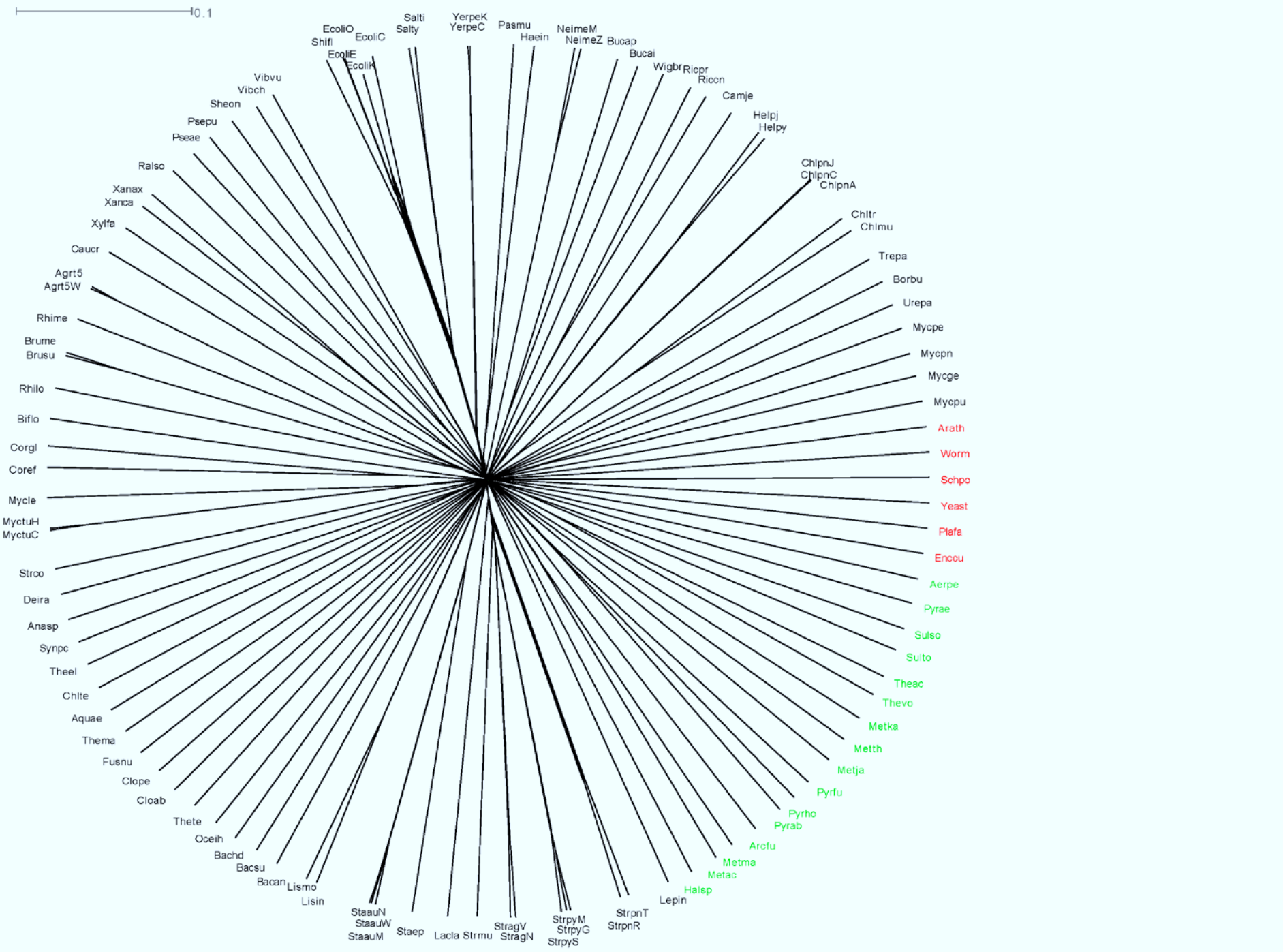


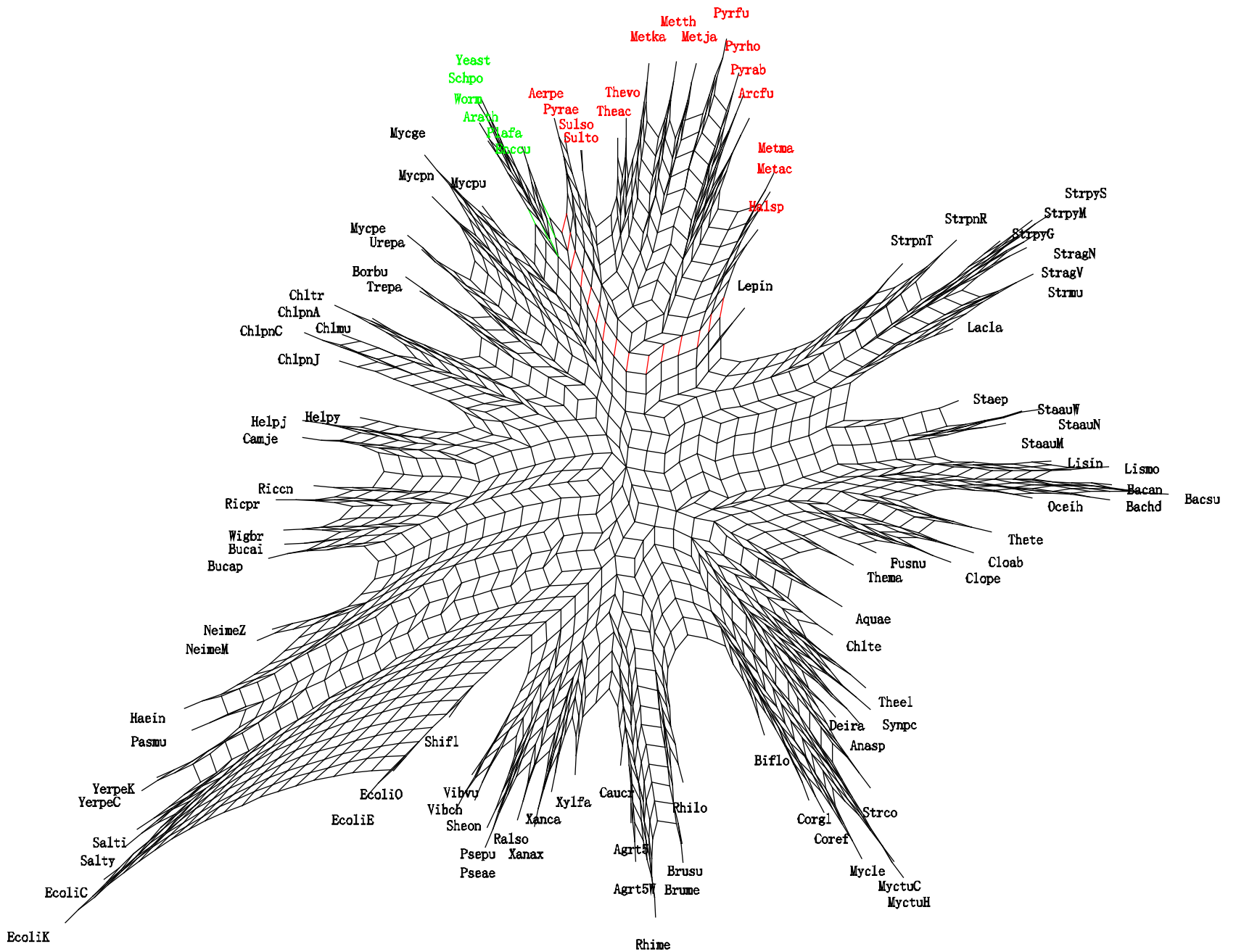


“Concentrating on the topology of the trees in the first place, we **did not scale** the branch lengths on the tree. However, these lengths should reflect evolution rates in terms of K -string composition changes. The **calibration** of branch lengths is further complicated by the overlapping nature of the K -strings when $K \geq 2$.”

The tree plotted to scale

0.1





- “However, the methods proposed so far are extremely crude. Oligonucleotide or oligopeptide frequencies are transformed into distances without any underlying model of evolution. It is nevertheless remarkable that something considered as a bias in standard sequence-based methods (Lockhart *et al.* 1994) contains a phylogenetic signal, but it is not yet clear whether accurate methods can be developed to extract it.”
(Philippe *et al.* 2005)
- “Why, for example, does the K -string complement of a proteome yield a tree that is similar to sequence-based trees?”
(Snel *et al.* 2005)

A heuristic probabilistic model

- Transition of frequency distribution:

Assume that the evolution process is stationary, and the average mutation rate per site is $1 - \alpha$. Let

$L = |\mathbf{S}| - K + 1$, then the conditional distribution

$$f\left(\mathbf{S}^{(t+1)}(\mathbf{i}) \mid \mathbf{S}^{(t)}(\mathbf{i}) = n\right) = \text{binomial}(n, \alpha^K)$$

$$* \text{binomial}\left(L, (1 - \alpha^K) \frac{\mathbf{E}\mathbf{S}(\mathbf{i})}{L}\right)$$

- Time evolution of expectation

$$\mathbb{E}\left(S^{(t')}(\mathbf{i}) \mid S^{(t)}(\mathbf{i}) = n\right) = \mathbb{E}S(\mathbf{i}) + \alpha^{K(t'-t)}(n - \mathbb{E}S(\mathbf{i}))$$

Calibration of CVTree

- Generalized CV

$$w(\mathbf{S}, \mathbf{i}) = (\mathbf{S}(\mathbf{i}) - \mathbf{S}^0(\mathbf{i}))c(\mathbf{i})$$

where $\mathbf{S}^0(\mathbf{i})$ is an unbiased estimator of $\mathbf{S}(\mathbf{i})$,
and $c(\mathbf{i})$ is arbitrary “normalization” factor.

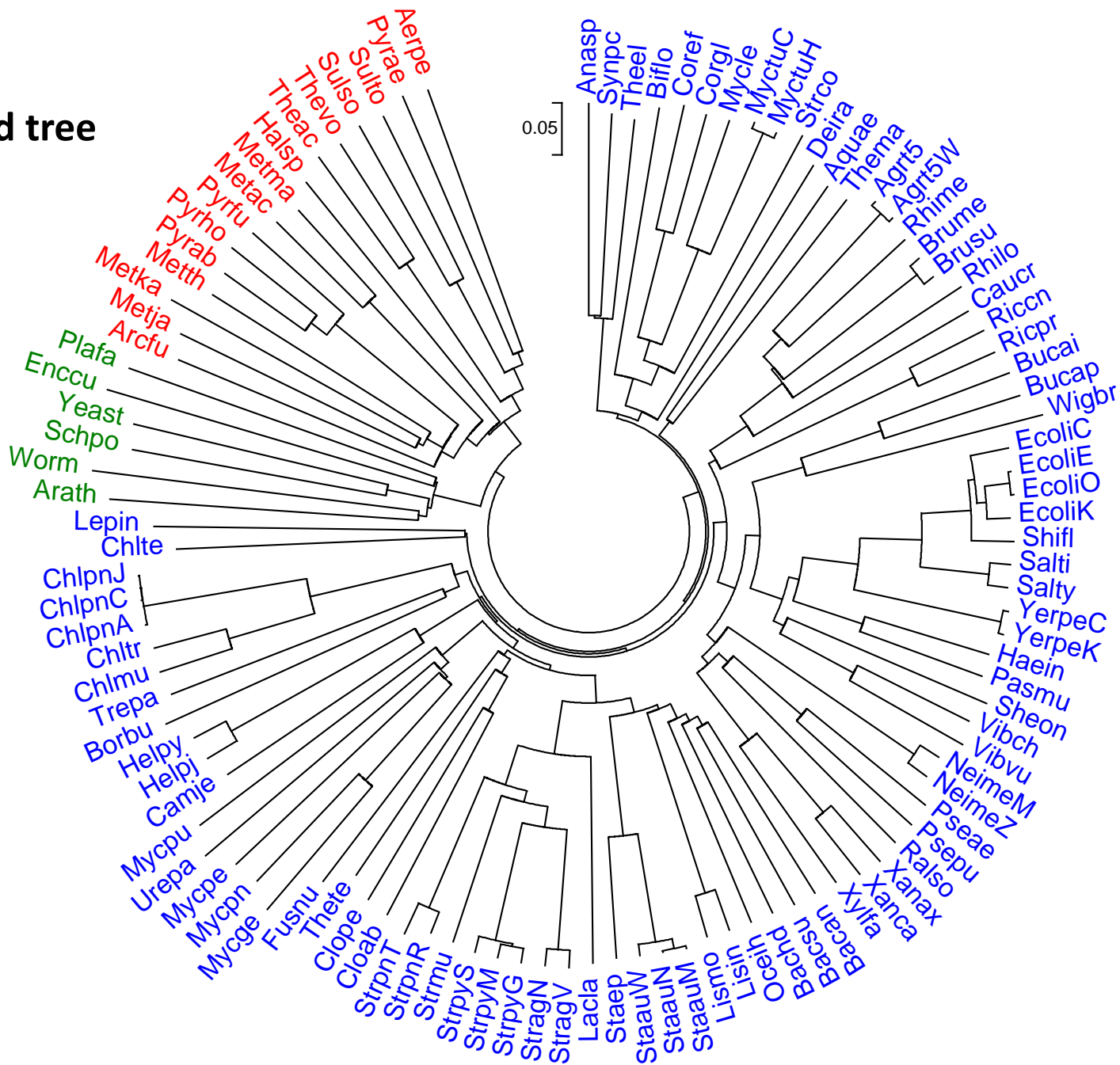
- Angle between vectors vs. sequence identity

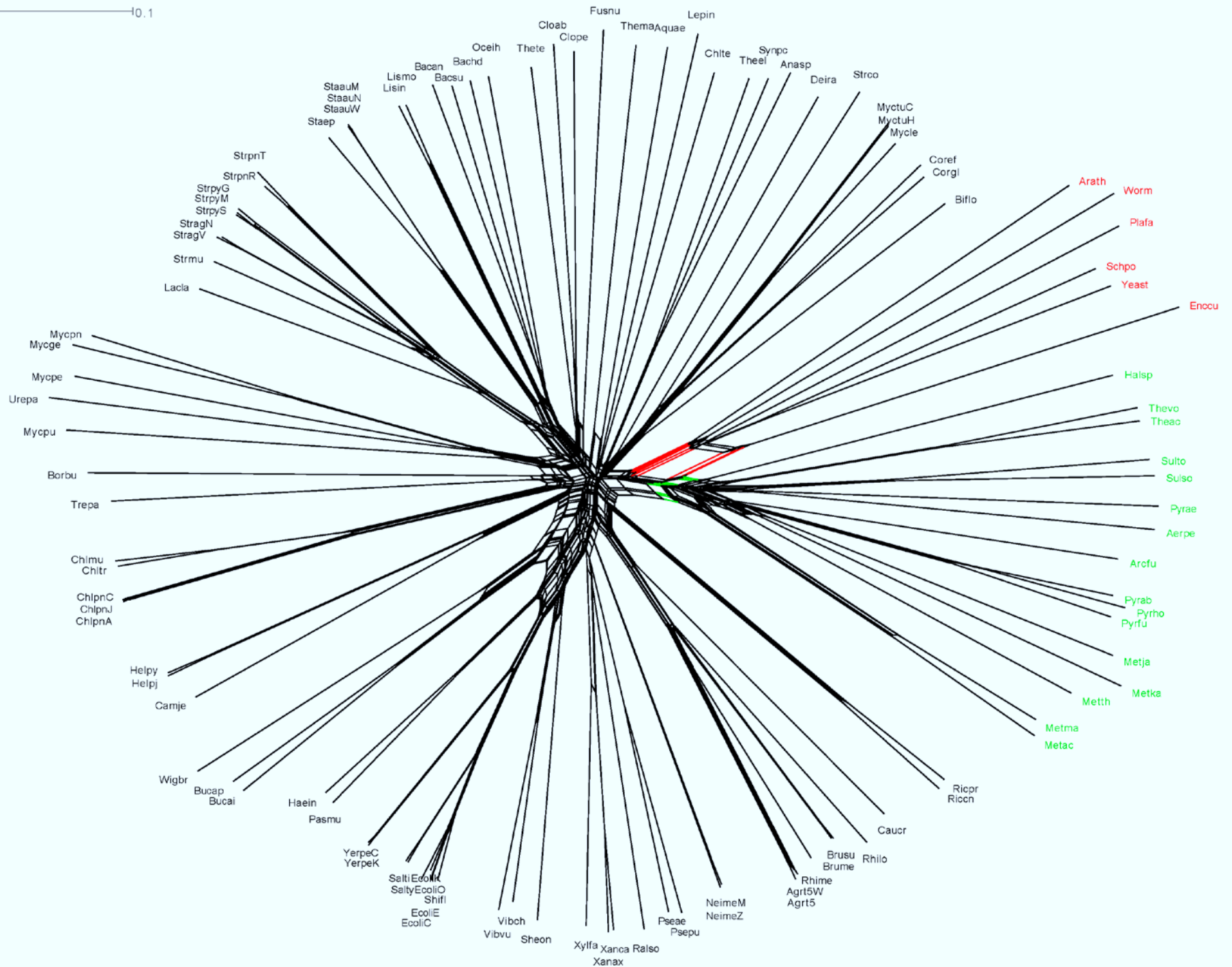
$$E \cos(w(\mathbf{S}), w(\mathbf{S}')) = q^K$$

- Calibration

$$\hat{p}(S, S') = 1 - \cos^{1/K}(w(S), w(S'))$$

The calibrated tree





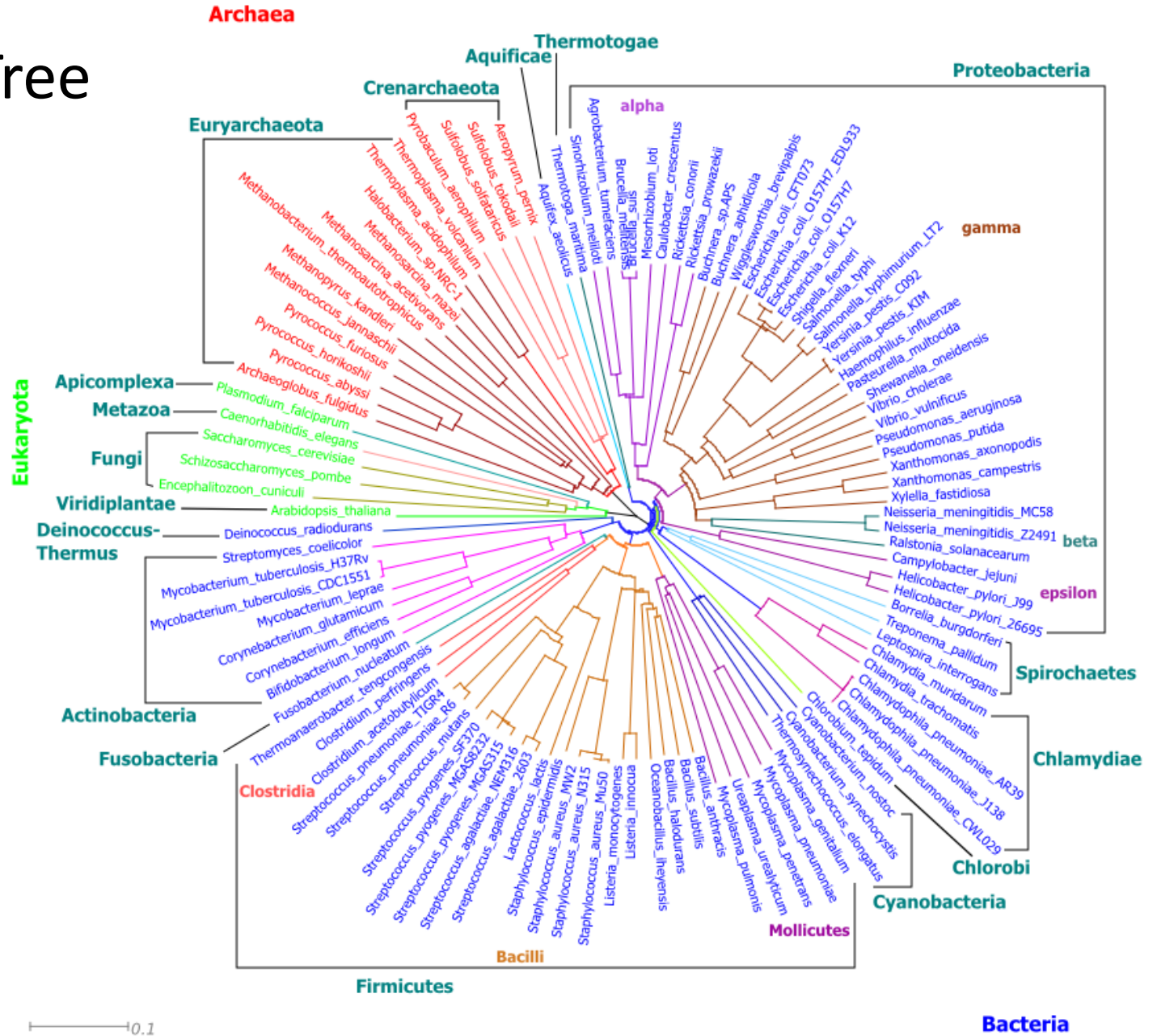
Using pre-/absence of words

- Any sequence S as above can be represented also by just the set $X(S)$ of its K -strings. In a “gedanken alignment”, ignoring repeats and collisions, we may put

$$Q^{(K)}(S, S') := \frac{1}{2} \left(\frac{|X(S) \cap X(S')|}{|X(S)|} + \frac{|X(S) \cap X(S')|}{|X(S')|} \right)$$

and $q^{(K)}(S, S') := Q^{(K)}(S, S')^{1/K}$

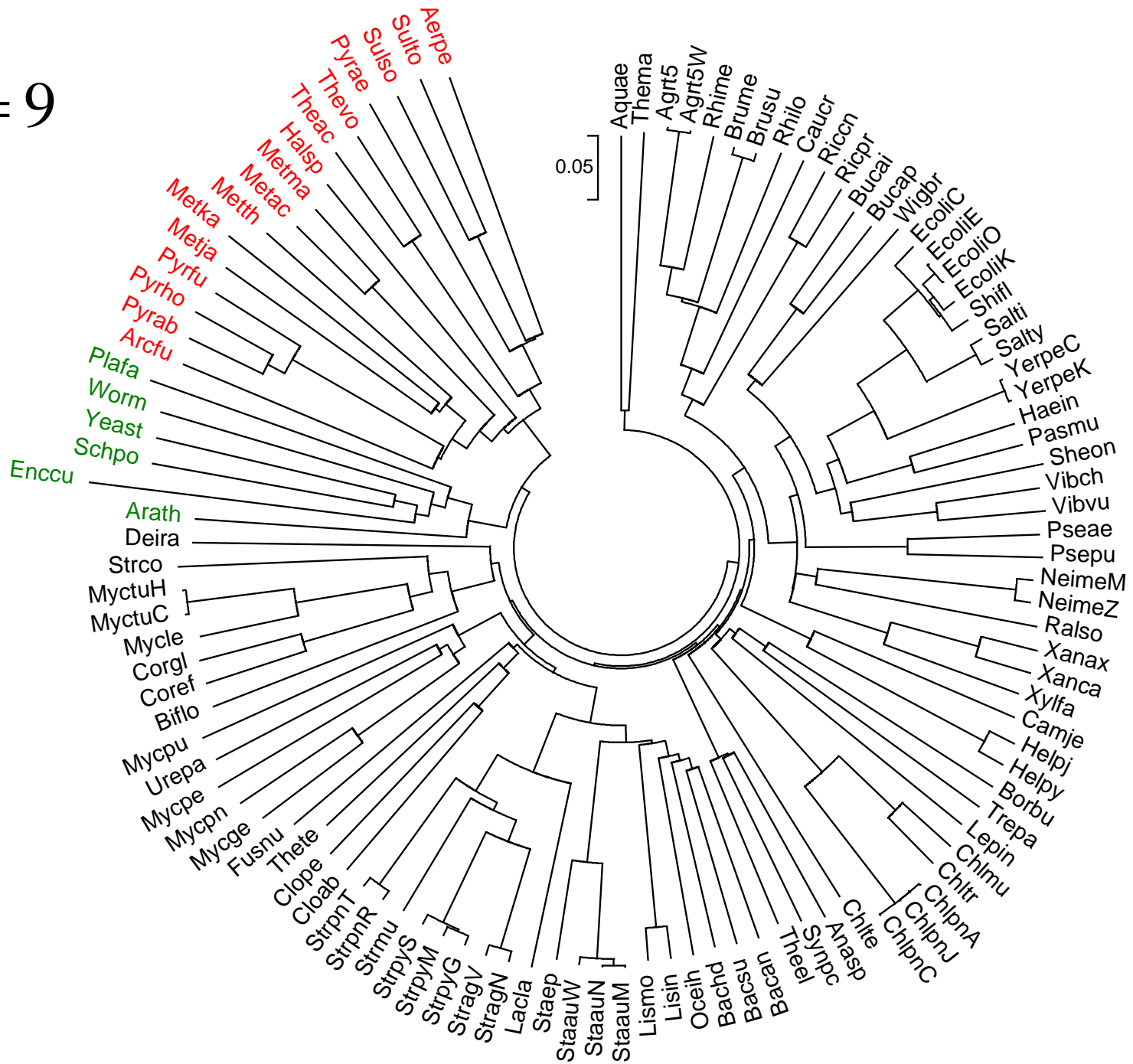
K = 9 Tree



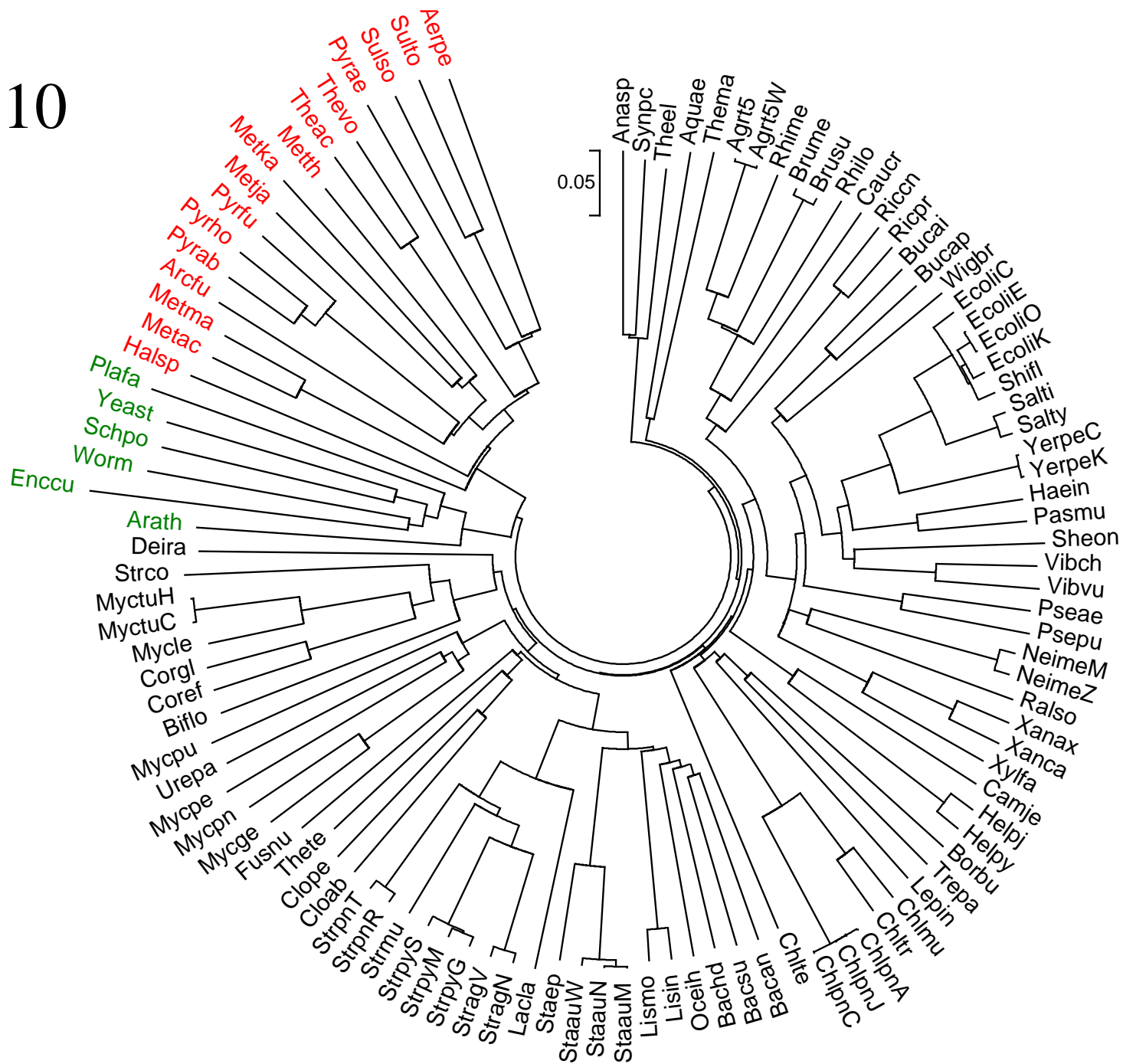
0.1

Bacteria

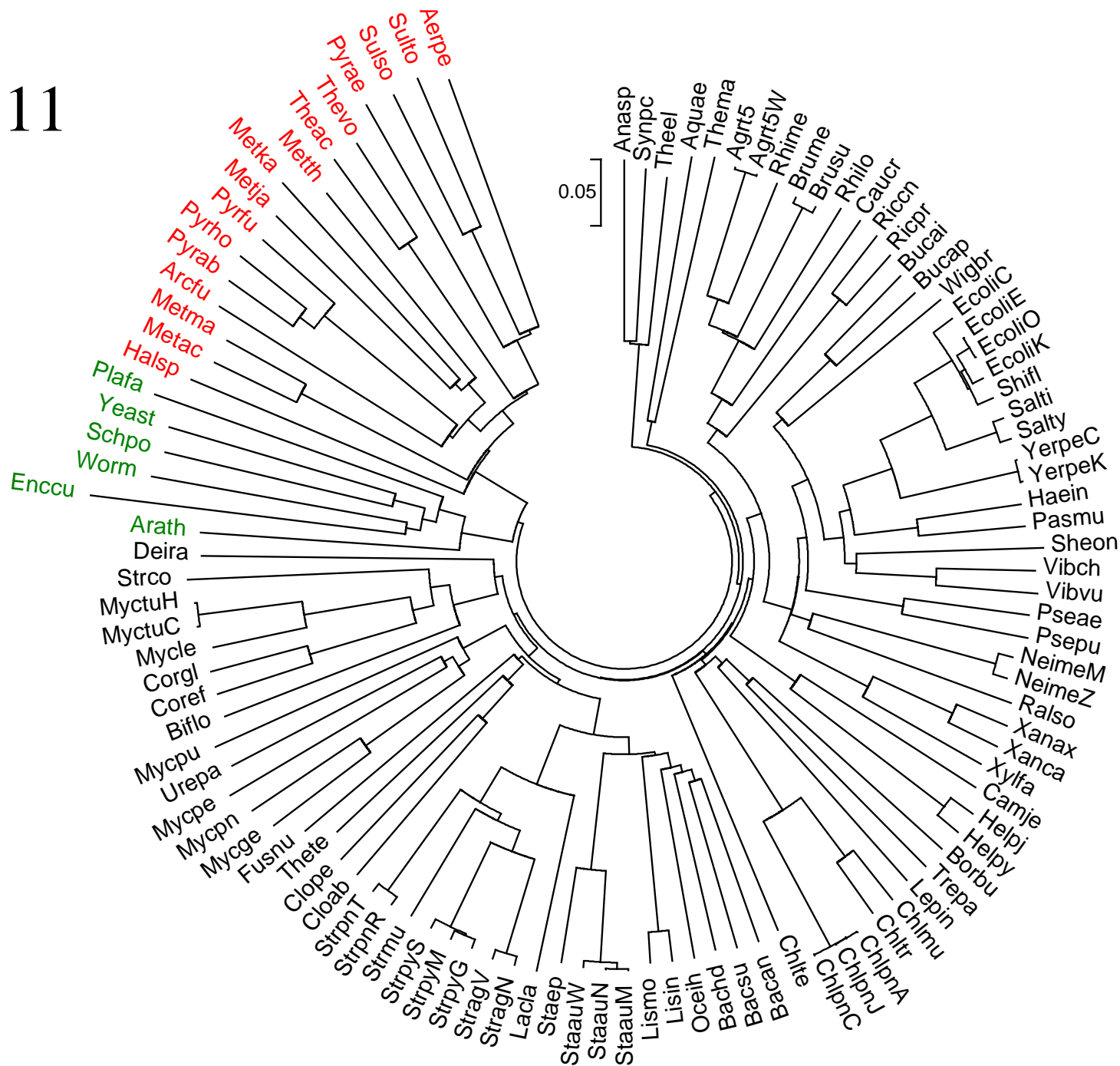
$K = 9$



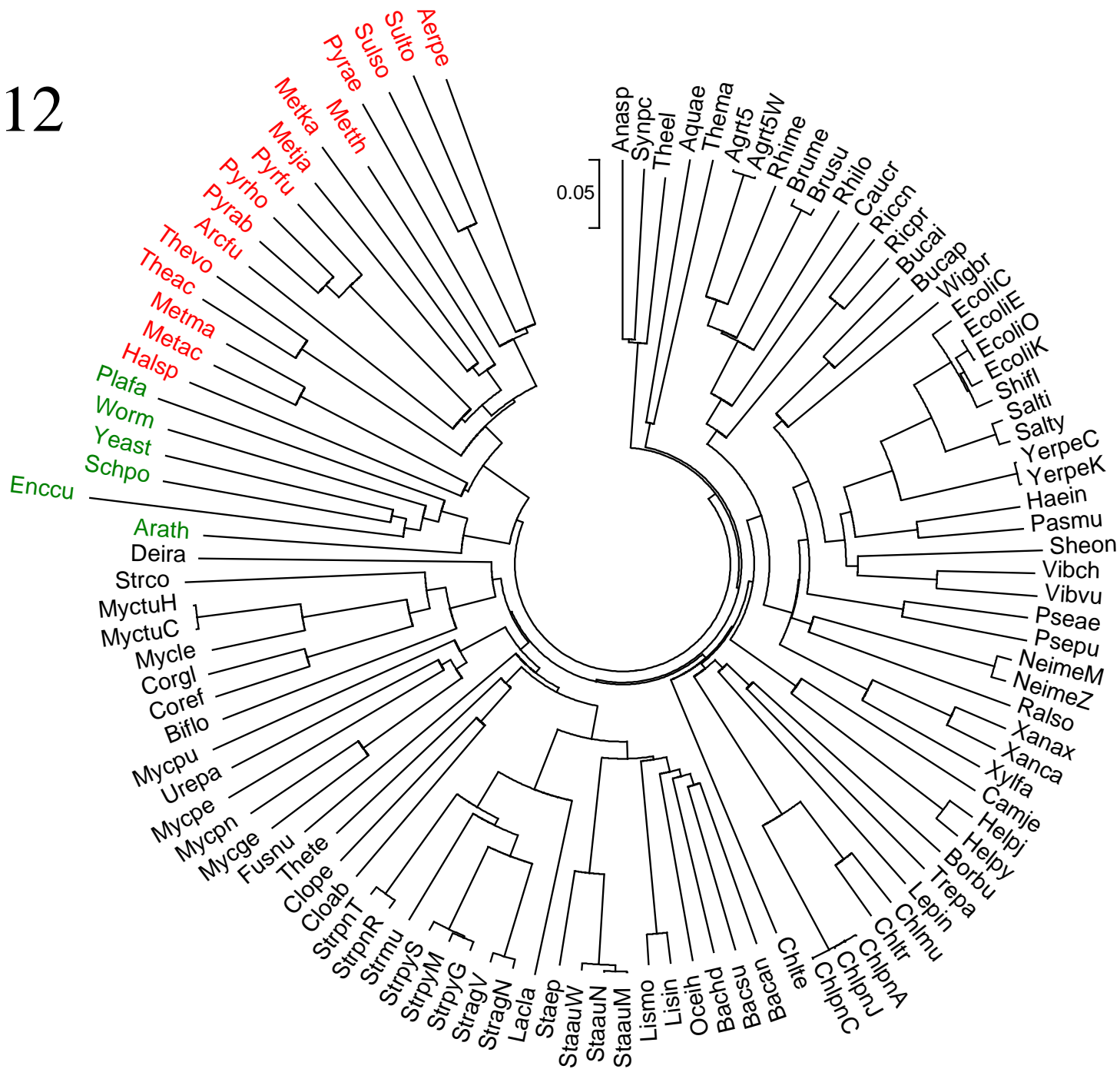
$K = 10$



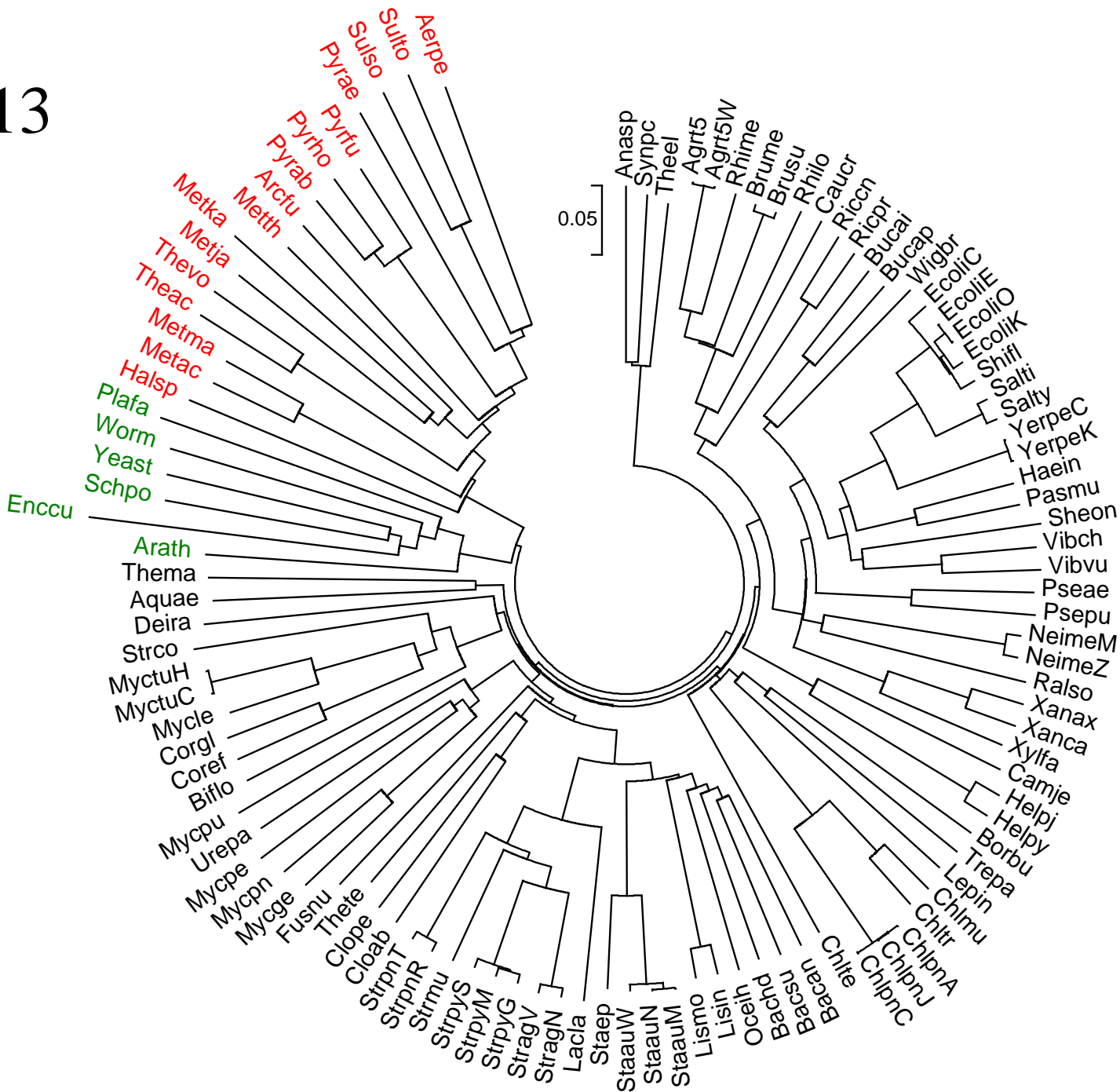
$K = 11$



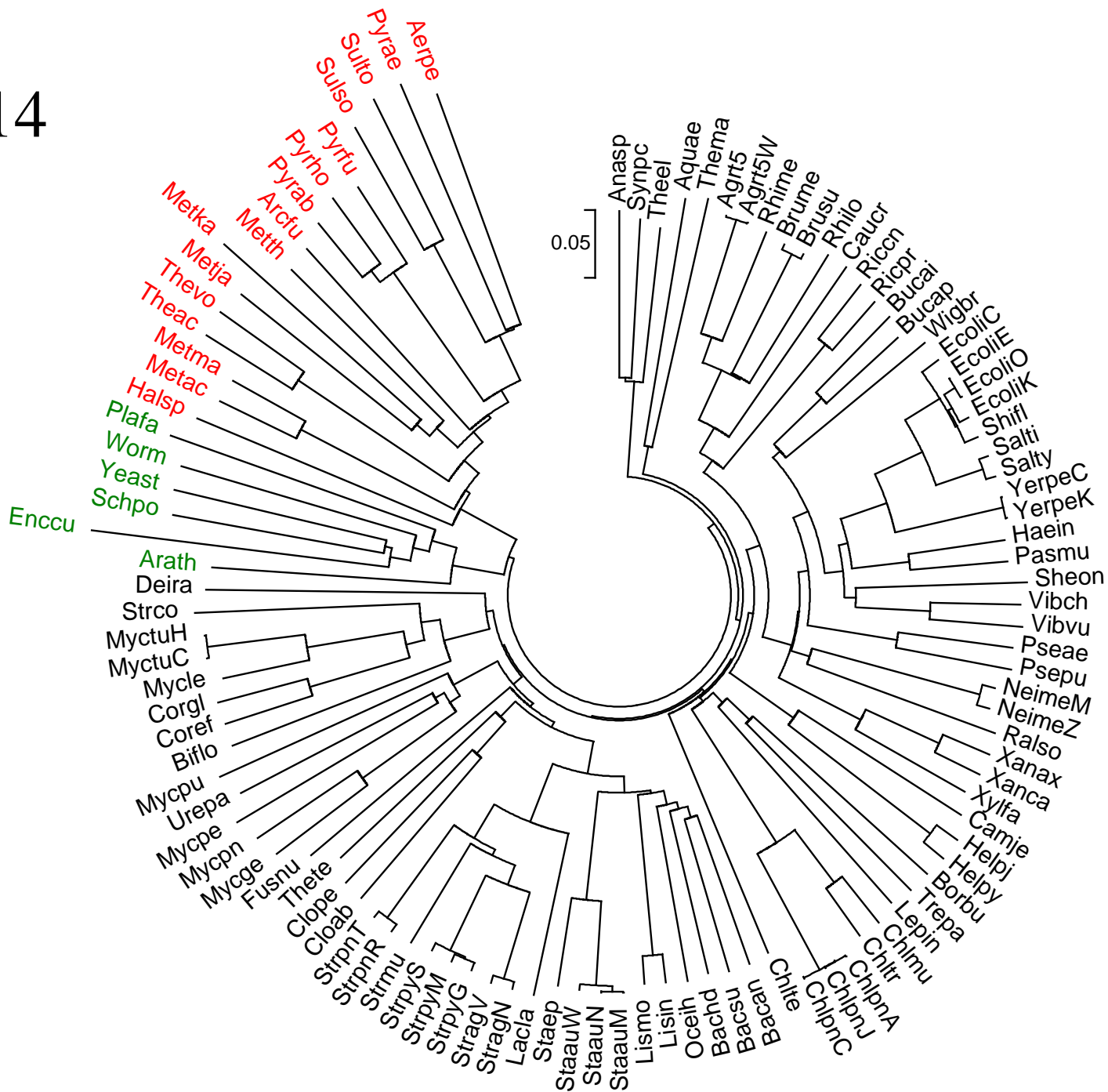
$K = 12$

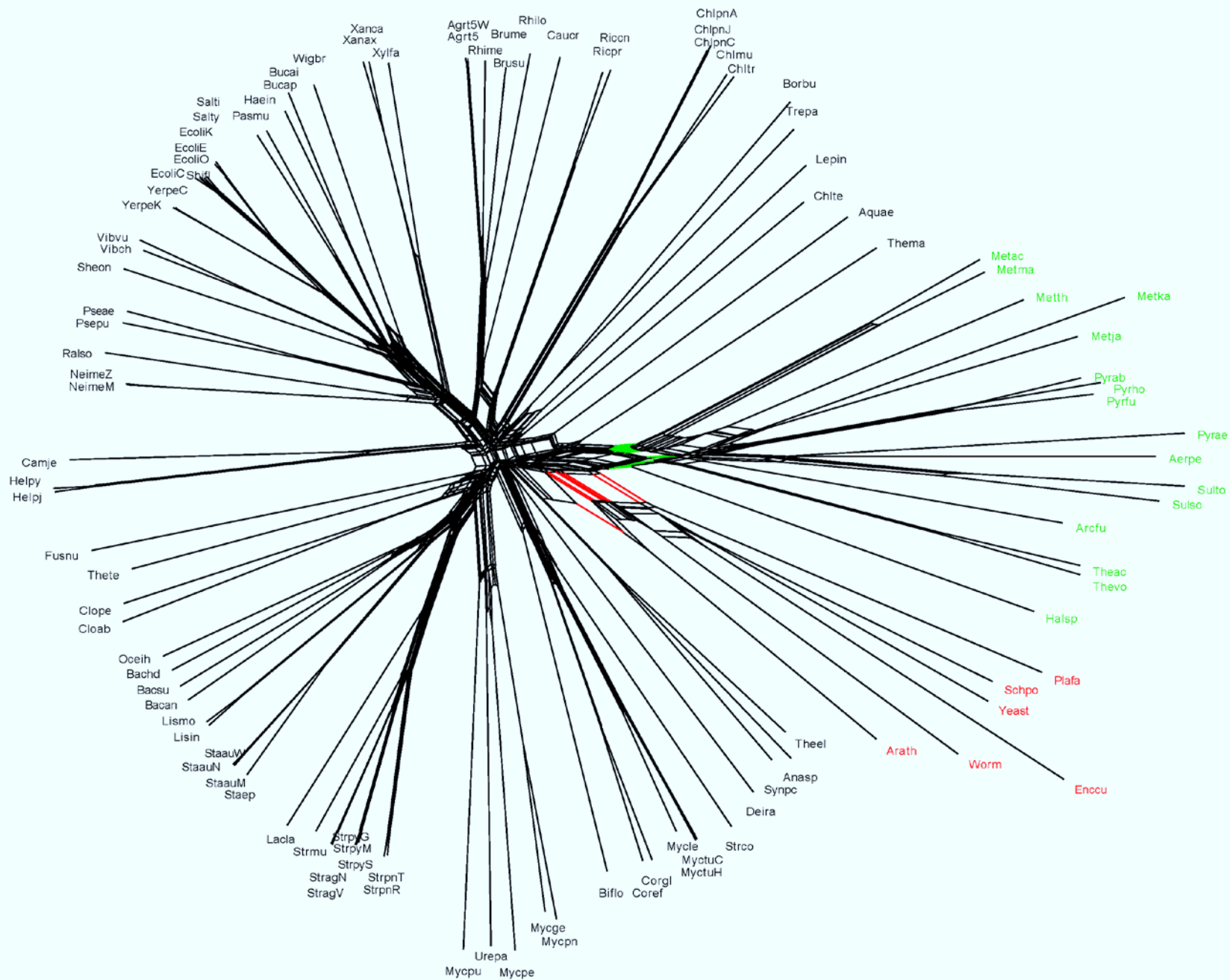


$K = 13$



$K = 14$





Using pre-/absence of words

- Any sequence S as above can be represented also by just the set $X(S)$ of its K -strings. In a “gedanken alignment”, ignoring repeats and collisions, we may put

$$Q^{(K)}(S, S') := \frac{1}{2} \left(\frac{|X(S) \cap X(S')|}{|X(S)|} + \frac{|X(S) \cap X(S')|}{|X(S')|} \right)$$

and $q^{(K)}(S, S') := Q^{(K)}(S, S')^{1/K}$

or $:= Q^{(K+1)}(S, S') / Q^{(K)}(S, S')$

Making up alignments

- Pattern:

1: key, must match

0: space, to watch

111110

- “Align” pairwise:

For keys occur in both sequences, concatenate letters in spaces.

X: VKAAWG, HAGEYG...

Y: VKAAWS, HAGEYG...

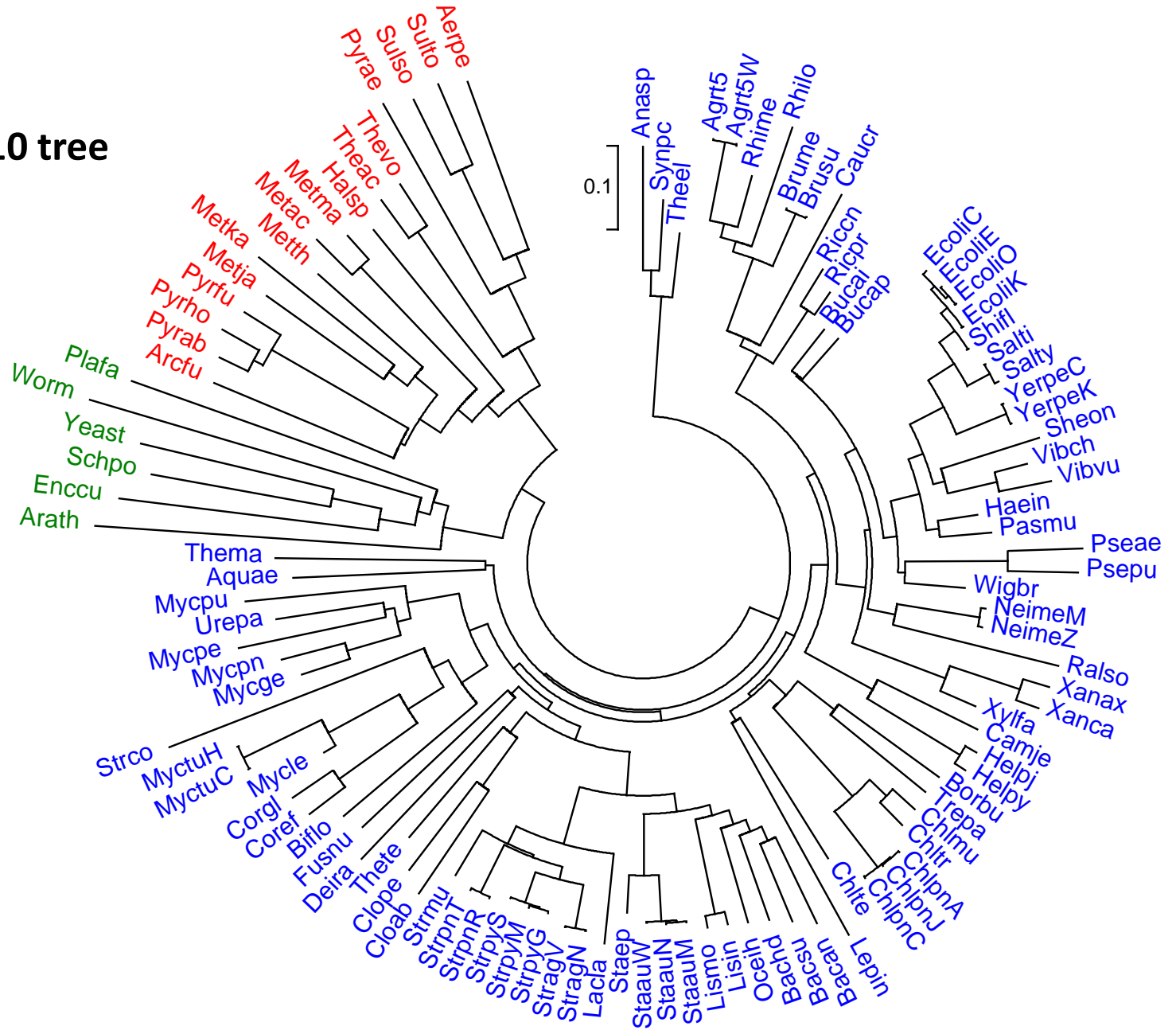
GG...

SG...

- Sequence-based distance

$$\hat{p} = n_d / n$$

011111110 tree



The missing words approach

- Probabilities of double- vs. single-missing

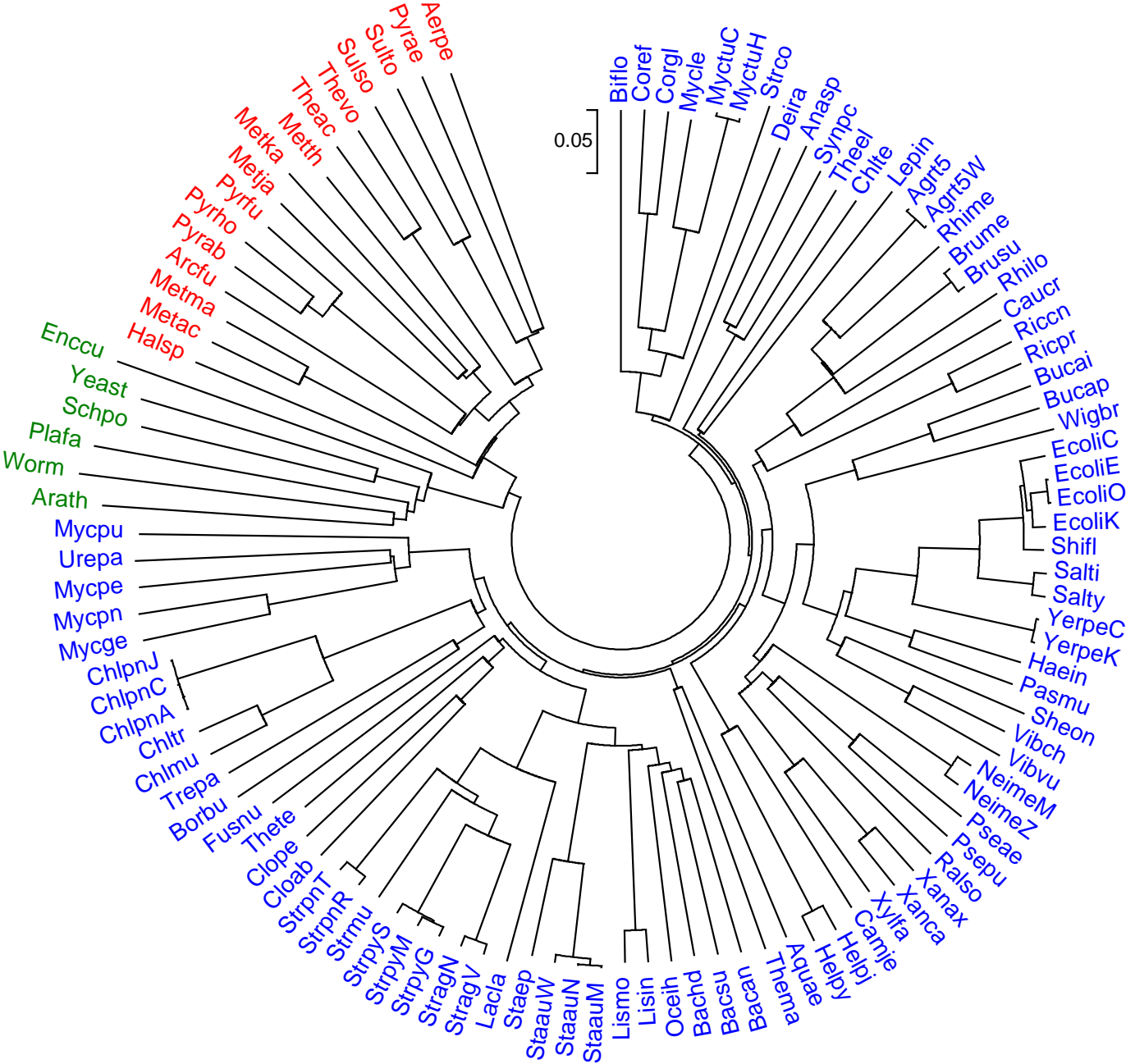
$$p_0 = u^{2-q^K} \Rightarrow \hat{q} = \left(2 - \frac{\ln p_0}{\ln u} \right)^{1/K}$$

- Estimation of probabilities

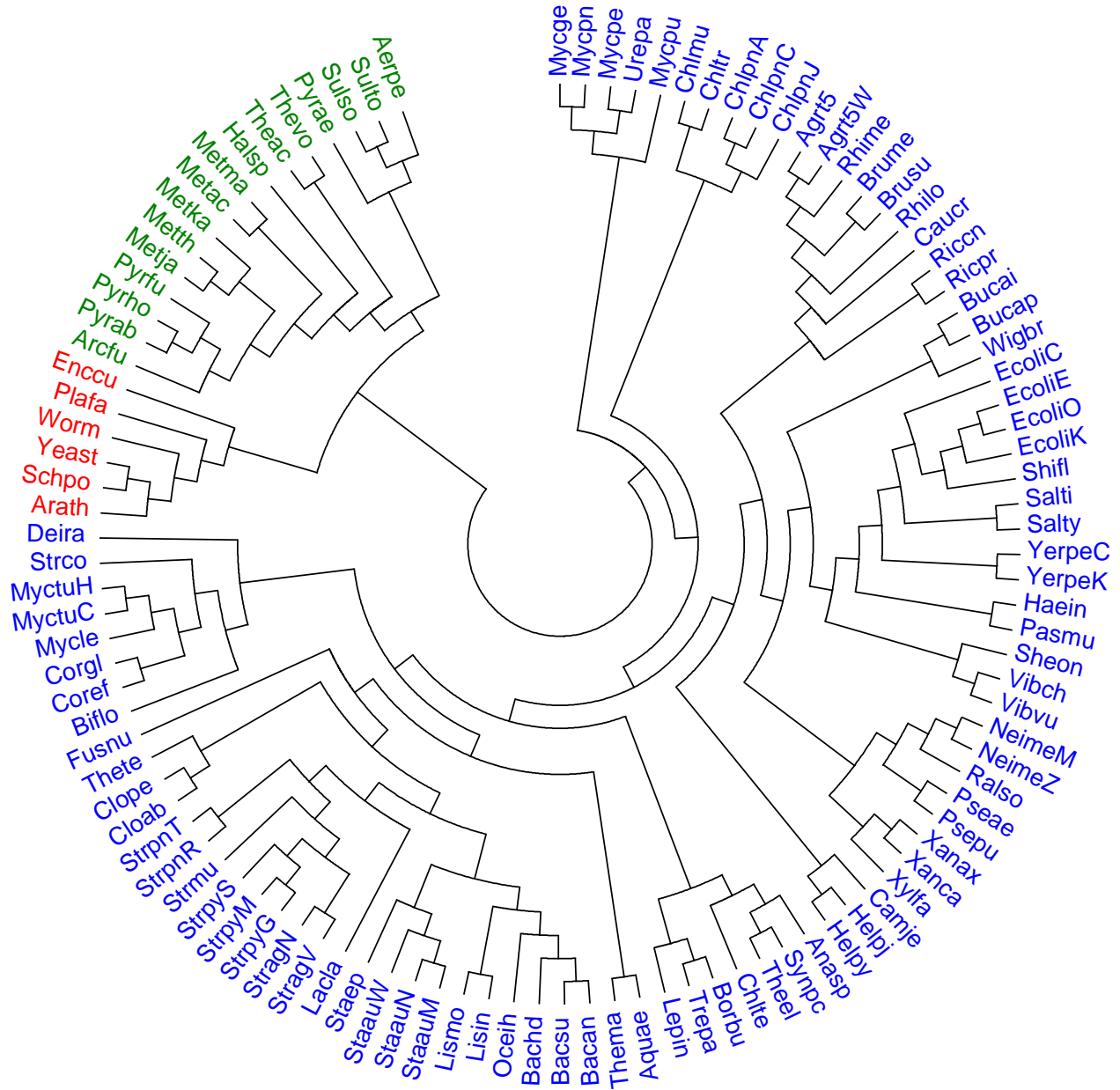
$$\hat{p}_0 = \frac{n_0}{|\Sigma|^K}$$

$$\hat{u} = \sqrt{\frac{n_0 + n_{\bar{B}}}{|\Sigma|^K} \cdot \frac{n_0 + n_{\bar{A}}}{|\Sigma|^K}}$$

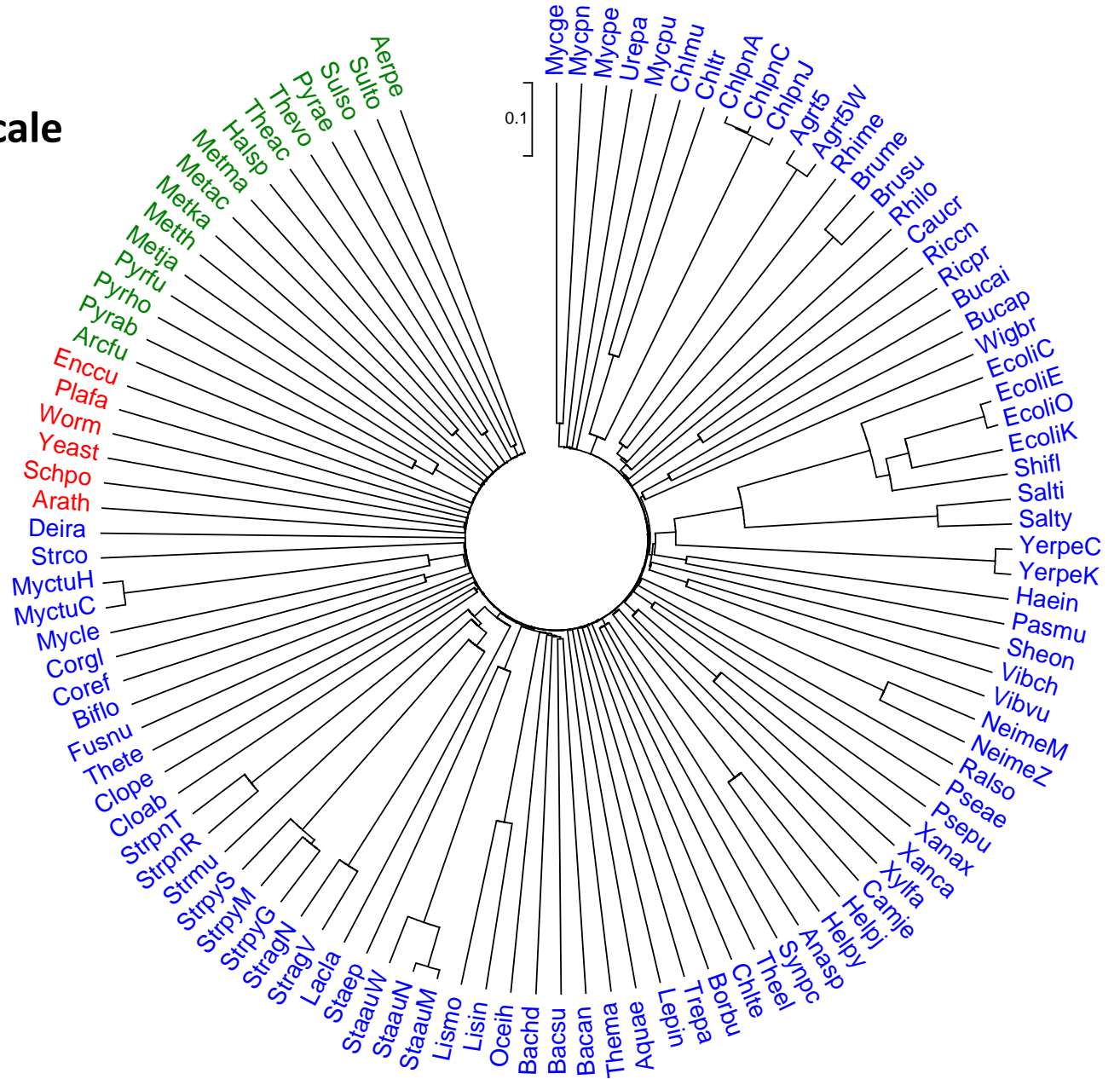
MW9 Tree



A tree based on
“Jensen-Shannon”
distance $JSh(S, S')$,
the “symmetrized”
Kullback-Leibler
distance between the
“induced K -string
probabilities”

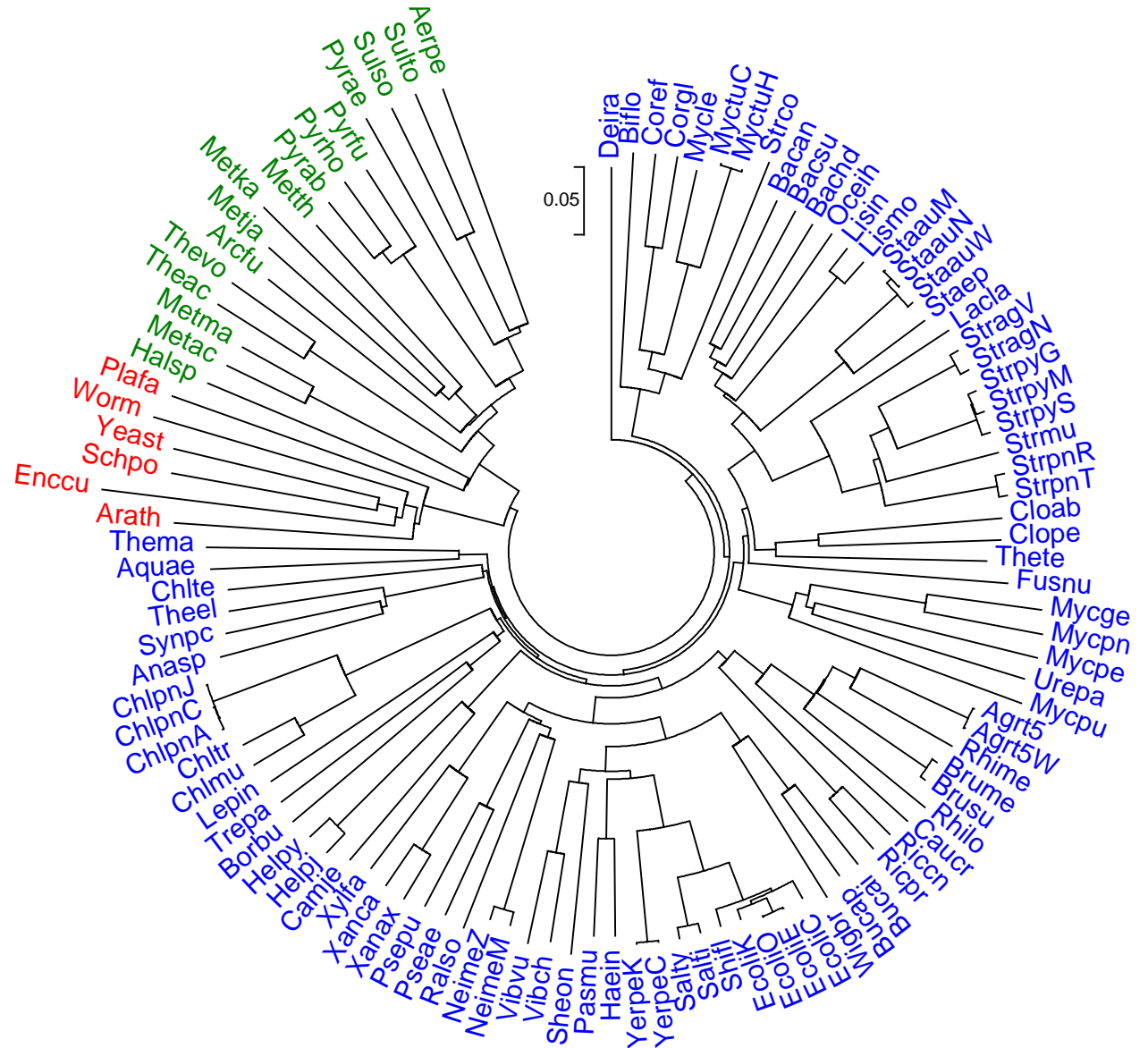


The tree drawn to scale

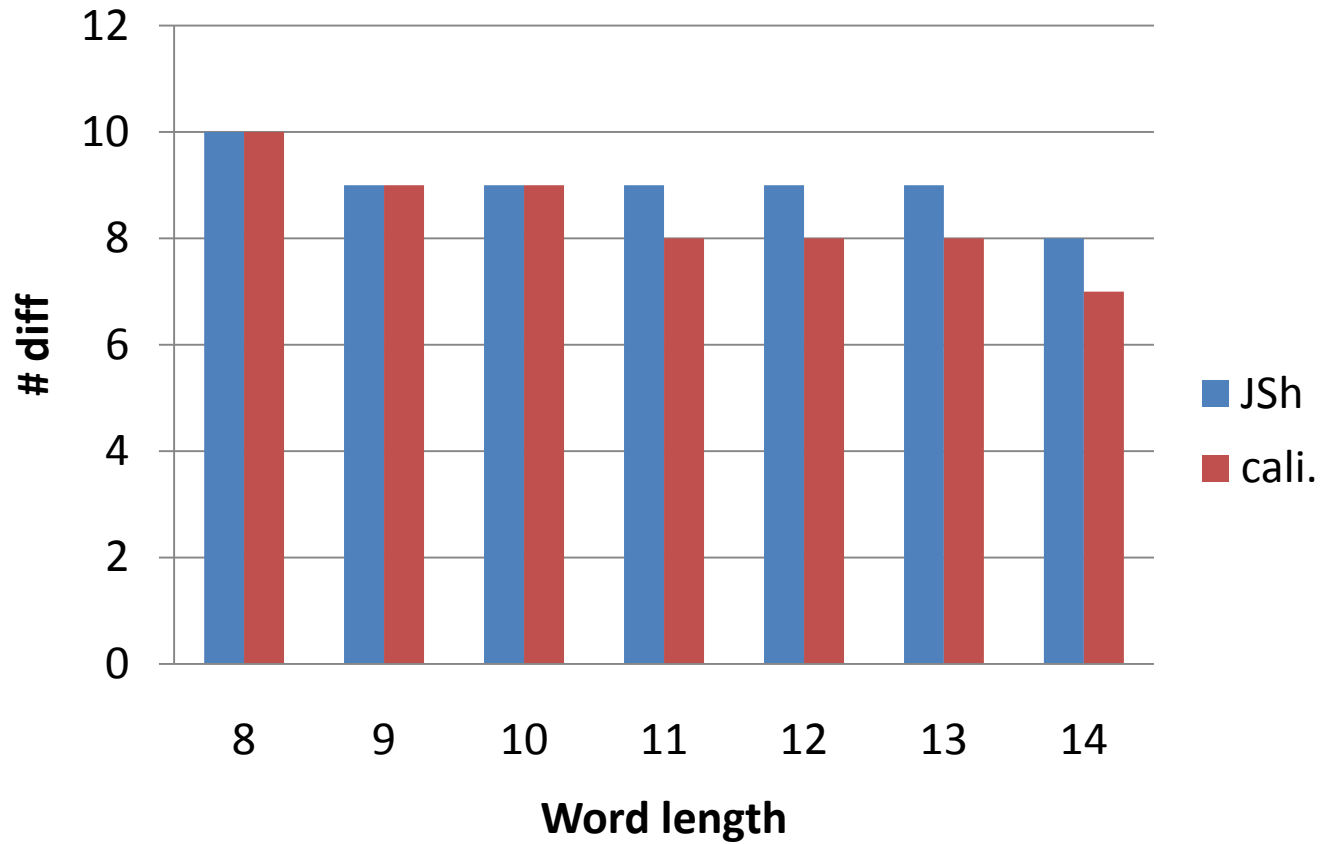


The tree using the calibrated distance

$$1 - (1 - JSh(S, S'))^{1/K}$$



Comparison with “standard” taxonomy



Outlook

- More realistic evolutionary dynamics
- Taking account of the variance of the frequency distribution
 - Better construction of composition vector
 - Optimal value of K

Acknowledgement

- Bailin Hao
- Andreas Dress
- Zhao Xu

Thank You!