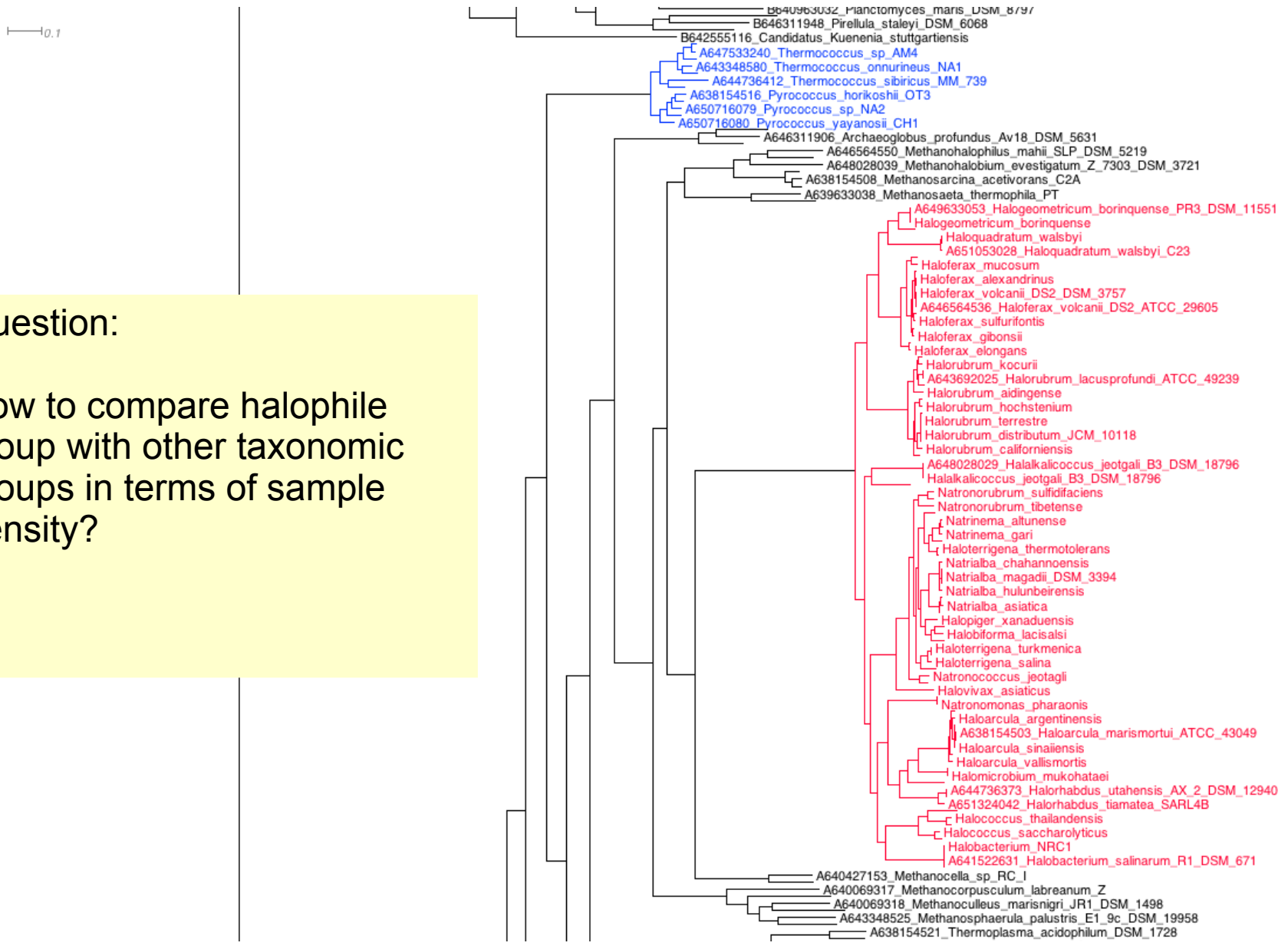


Phylogenetic Tree Branch Density Estimation & Tree-base OTU Grouping Revisit

Dongying Wu

July, 2012

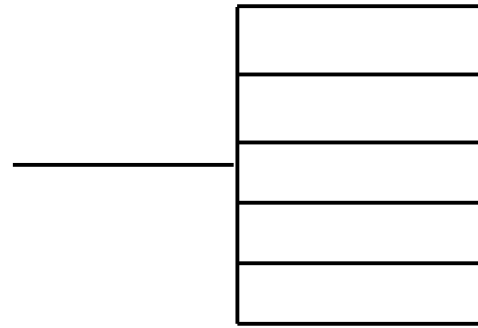
Fasttree of halophiles and all IMG bacteria and archaea based on the concatenated alignment of 40 markers



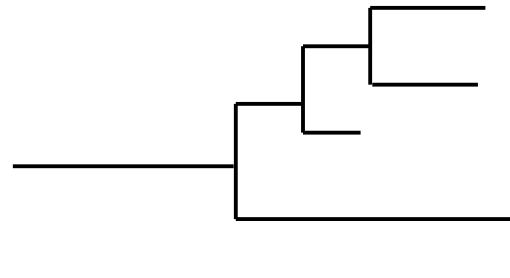
Question:
How to compare halophile group with other taxonomic groups in terms of sample density?

Estimate taxa density of a node in a tree

Taxa number as Density

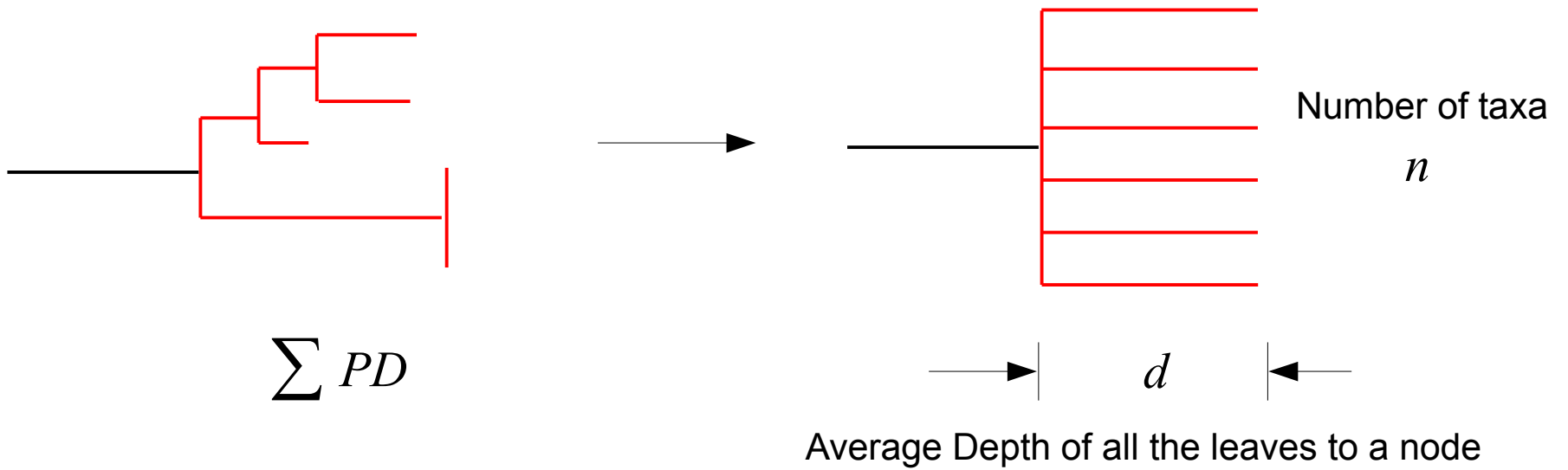


The underline assumption for taxa count as density



Density Calculation should take tree structure into account

Estimate taxa density of a node in a tree



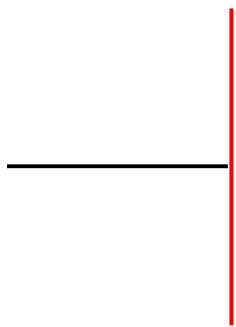
$$Density = \frac{\sum PD}{n \times d} \times n$$



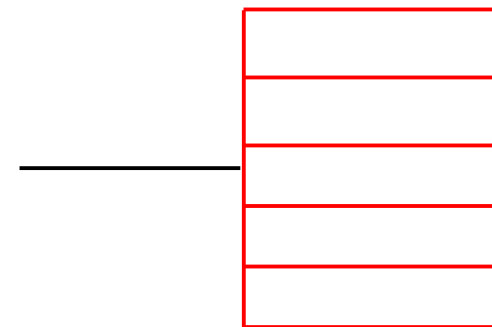
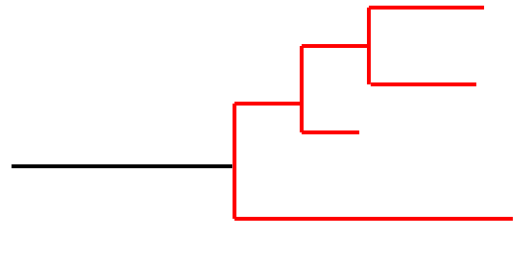
$$Density = \frac{\sum PD}{d}$$

$$Density = \frac{\sum PD}{d}$$

For a node with n leaves
Density is a number between 1 and n

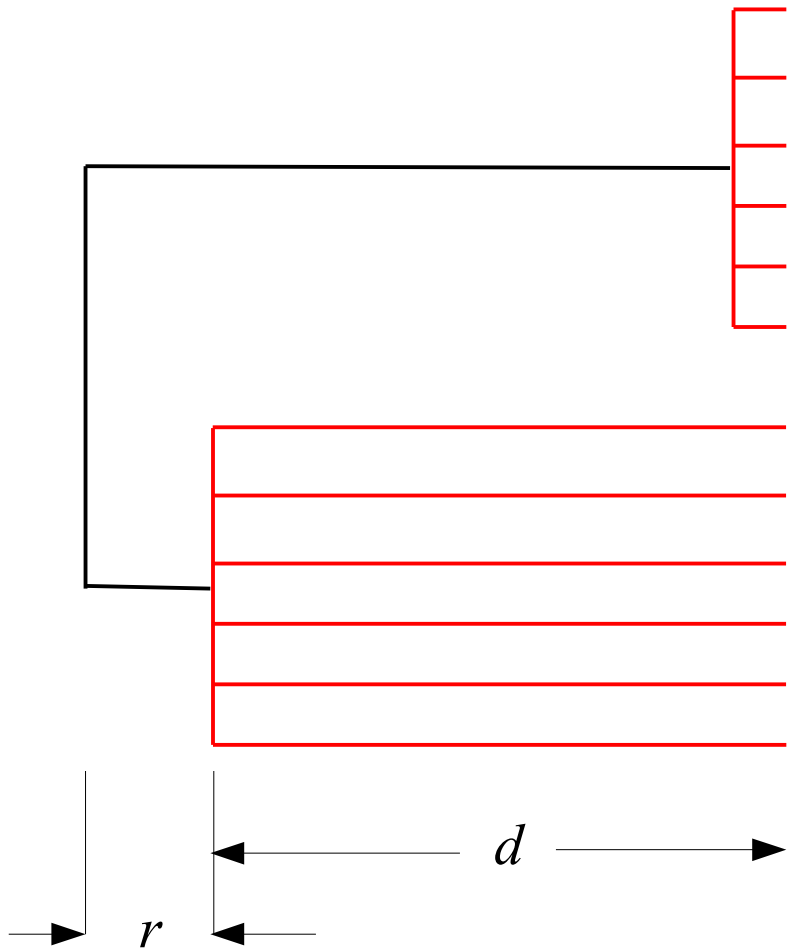


Density of 1



Density of n



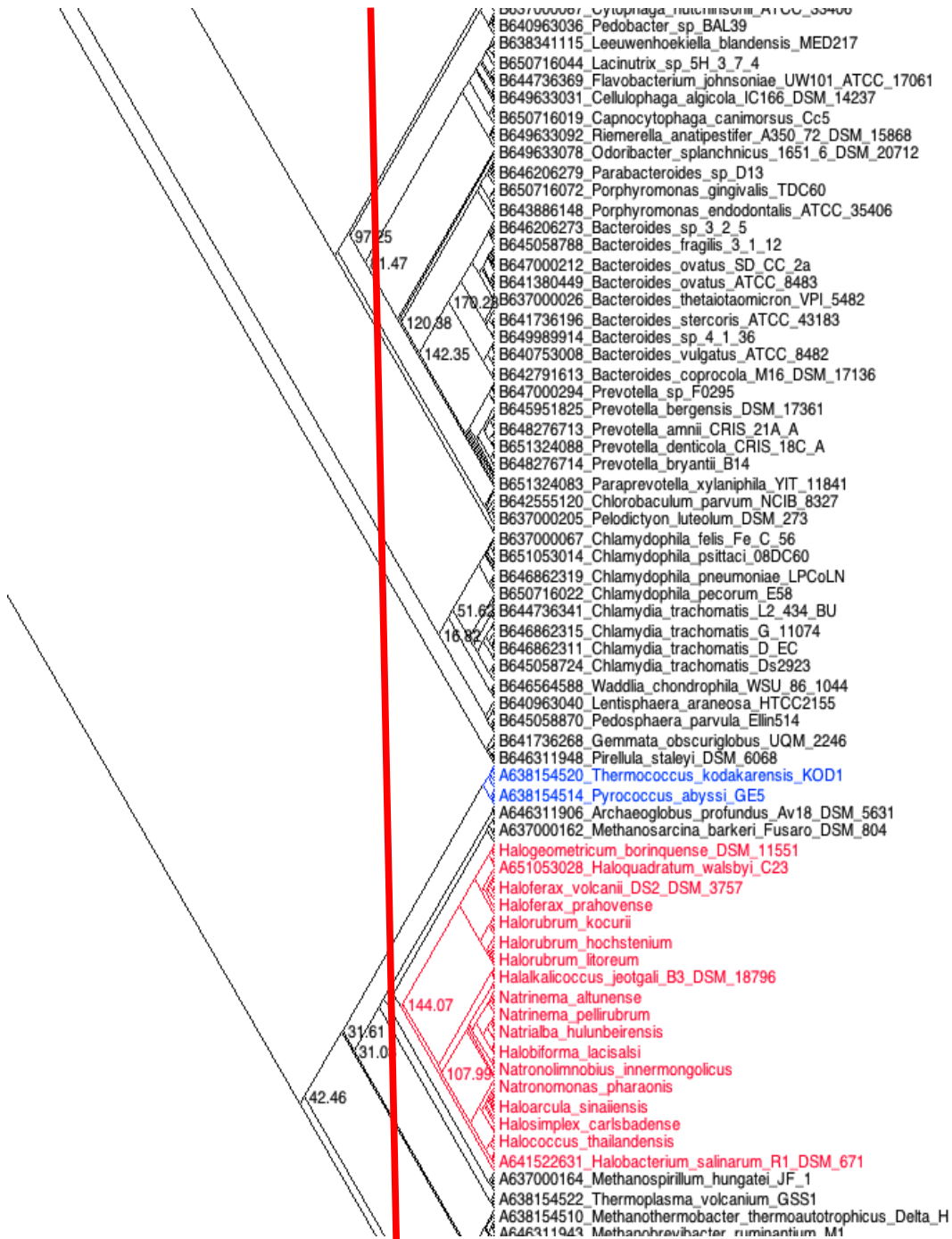


Density calculation need to take the relationship of root, node and leaves into account

$$Density = \frac{\sum PD}{d}$$

$$Relative\ Density = \frac{Density}{d/(r+d)}$$

$$Relative\ Density = \frac{(d+r) \times \sum PD}{d^2}$$



Halophile Group Relative Density: 144.07
 Thermococcus/Pyrococcus Group: 69.79

Does density comparison between nodes make biological sense?

Tree-base OTU grouping revisit

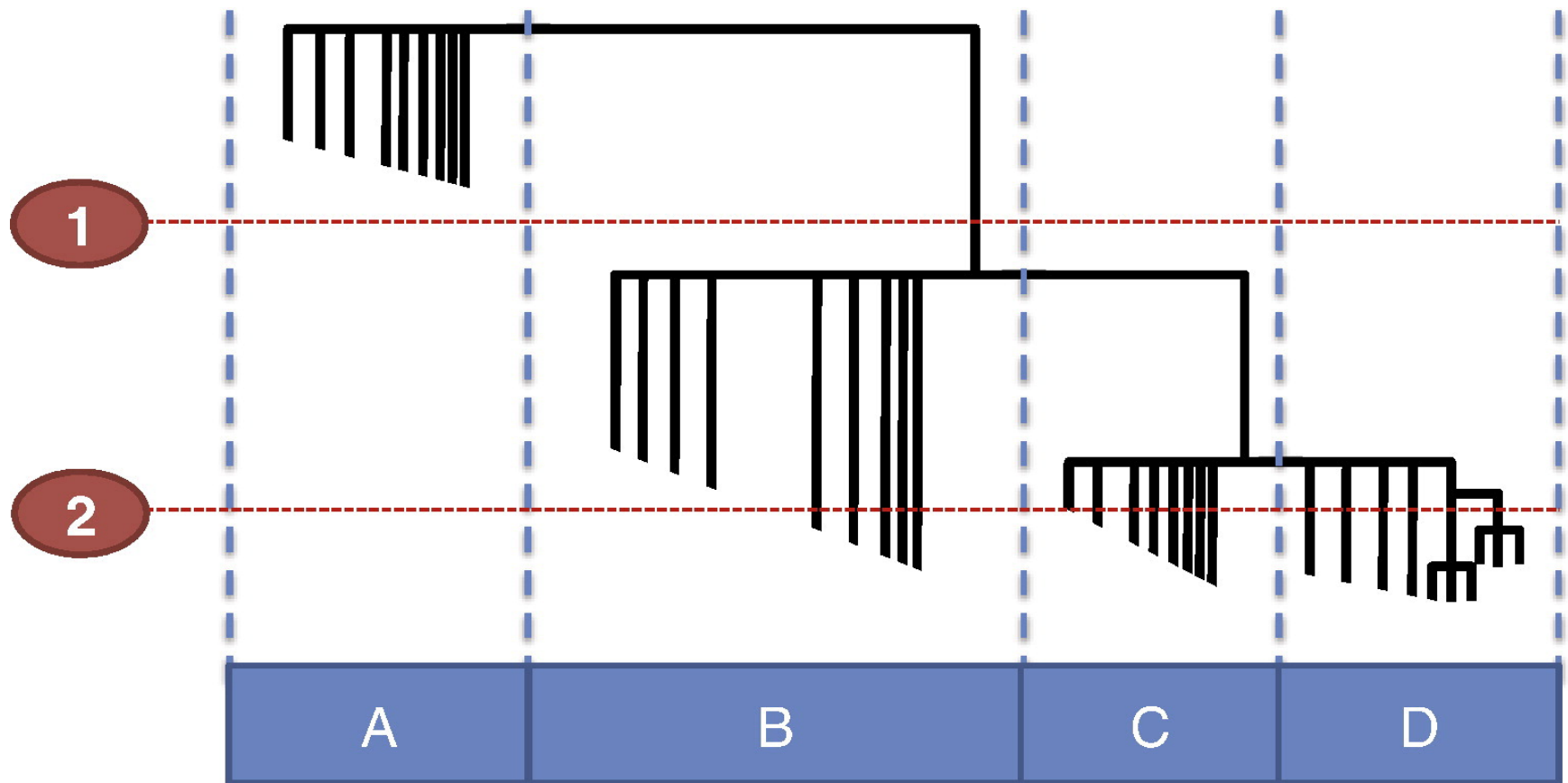
Version 1: For each node, calculate the average distance of all the leaves and the standard deviation, and decide if a node defines an OTU

Problem: Slow, RAM demanding

Outside Progress:

1. Tree topology based OTU grouping (Dynamic Tree Cutting, D. Serre)
2. Tree branch-length based OTU grouping (PHYLOSEQ, P. McMurdie & S. Holmes)

Dynamic Tree Cutting



Genomics. 2011 Oct;98(4):253-9. Epub 2011 Apr 15.

A novel method for determining microflora composition using dynamic phylogenetic analysis of 16S ribosomal RNA deep sequencing data.

Chan ER, Hester J, Kalady M, Xiao H, Li X, Serre D.

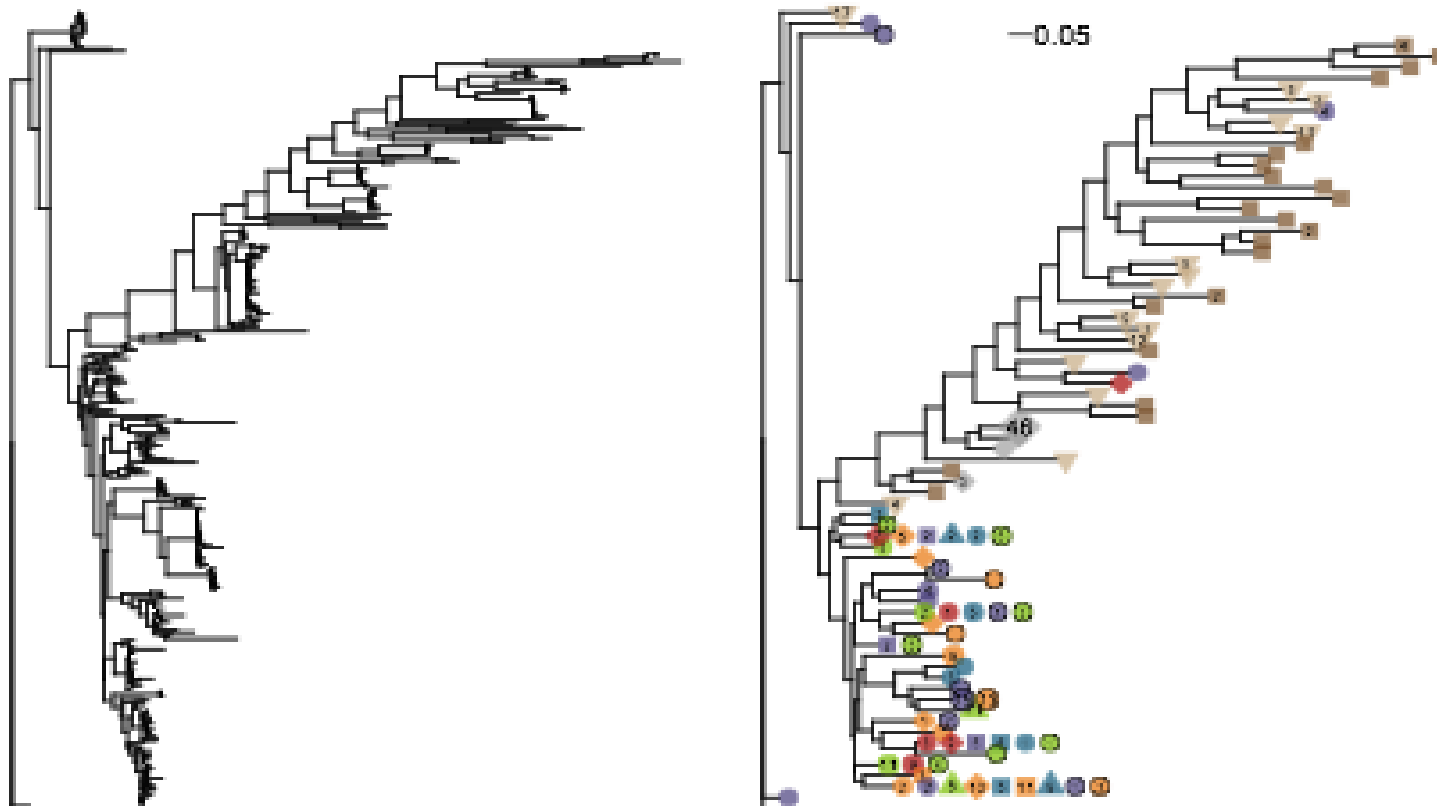
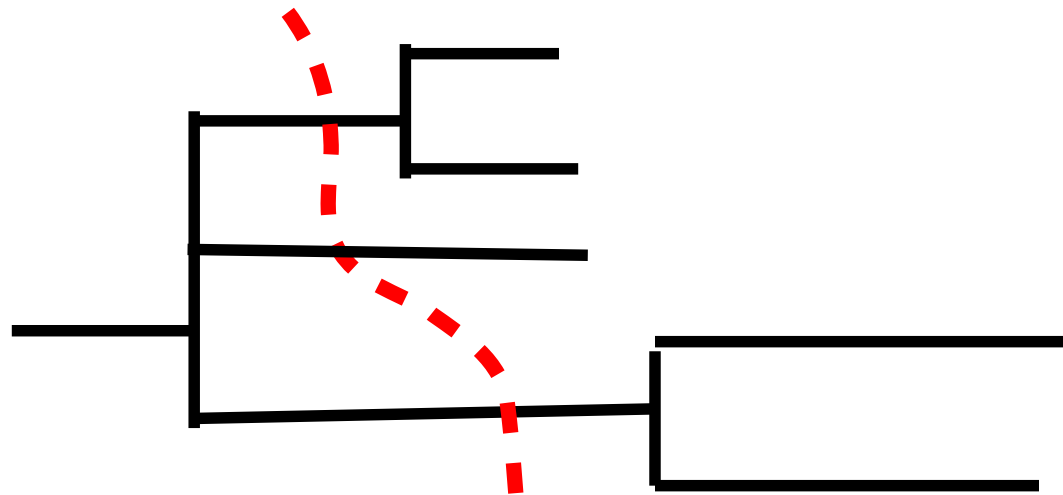


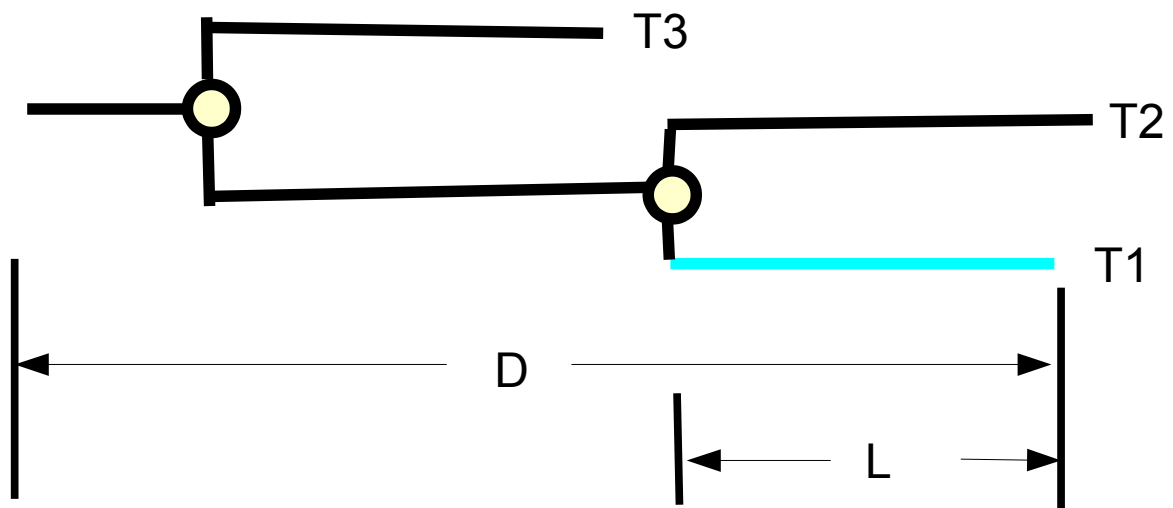
Fig. 6: Example of phylogenetic sequence data before and after basic clustering with `tipglom()` function. (Left) Standard phylogram produced using default plotting function and no OTU clustering. (Right) Annotated phylogram after OTU clustering with `tipglom()`. Different symbols next to each tip indicate different samples in which the OTU was observed. The number inside each symbol indicates the respective number of individuals of a given OTU were observed in each sample.

Tree-based OTU classification

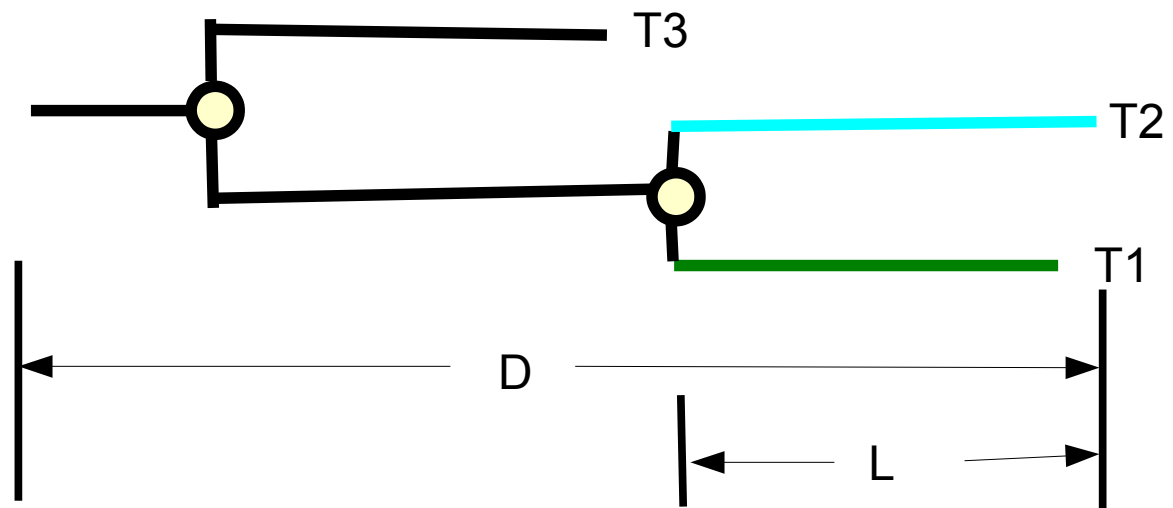


Rule 1: The classification should be based upon branch length rather than tree topology

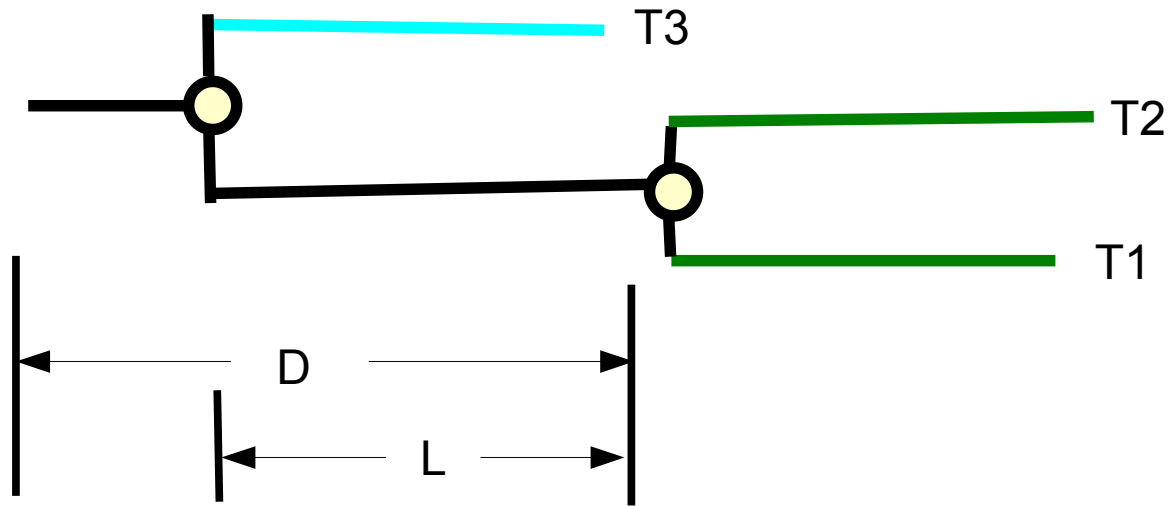
Rule 2: Dynamic cutting (branch length should be normalized to deal with different evolution rates from different lineages)



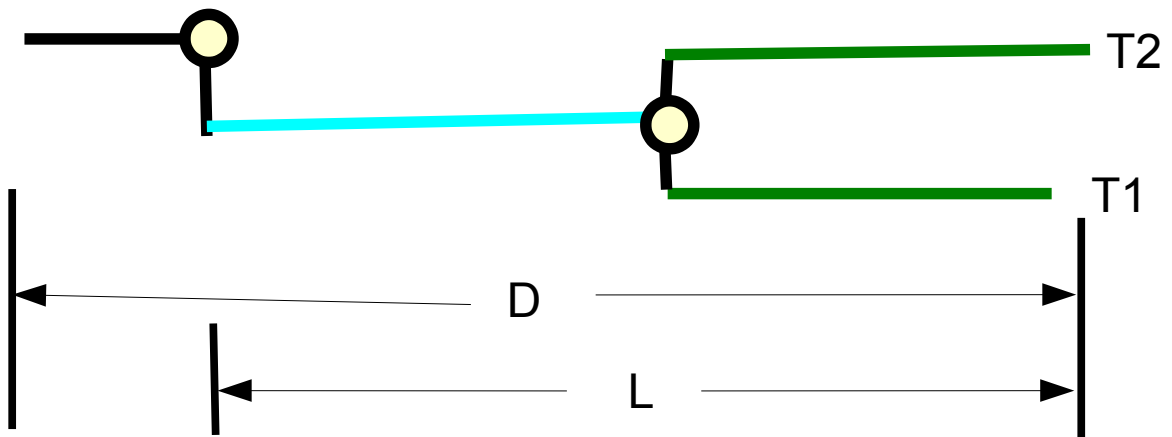
$L/D \geq \text{cutoff} ?$ No



$L/D \geq \text{cutoff} ?$ No

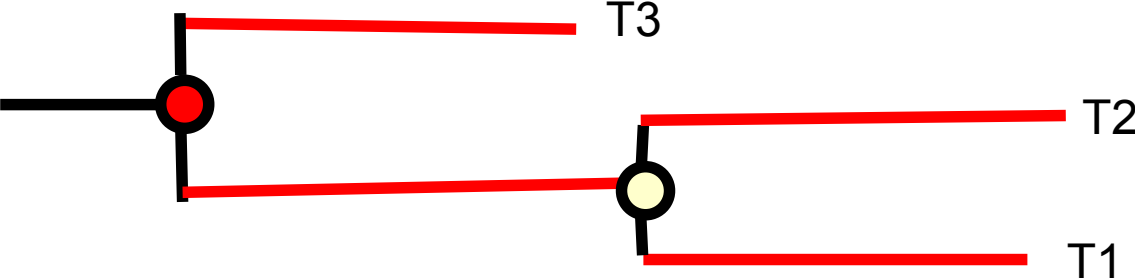


$L/D \geq \text{cutoff} ?$ Yes \longrightarrow (T3)



$L/D \geq \text{cutoff} ?$ Yes \longrightarrow (T1,T2)

How to calculate the average depth from tips to a node?

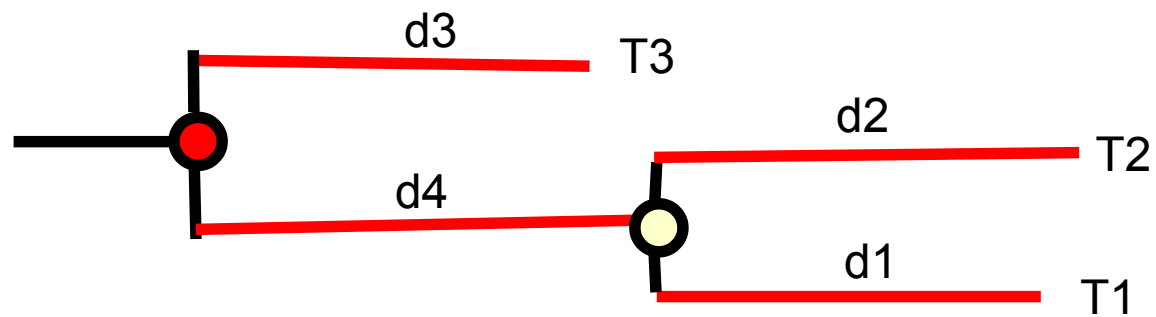


Naive Average:



Weighted Average:

Weight factor: The true contribution of a taxon to a node



T3 weight factor : 1

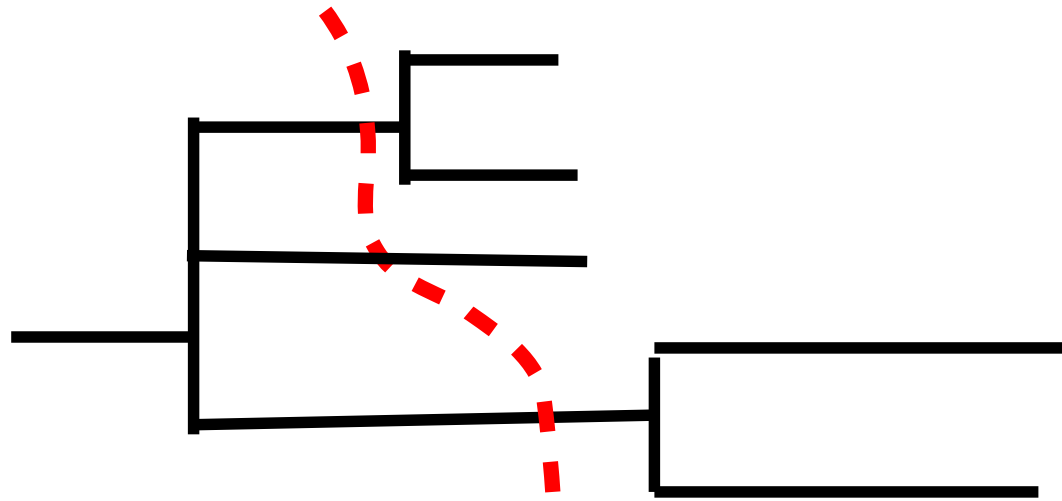
T2 weight factor : $(d2+d4/2)/(d2+d4)$

T1 weight factor : $(d1+d4/2)/(d1+d4)$

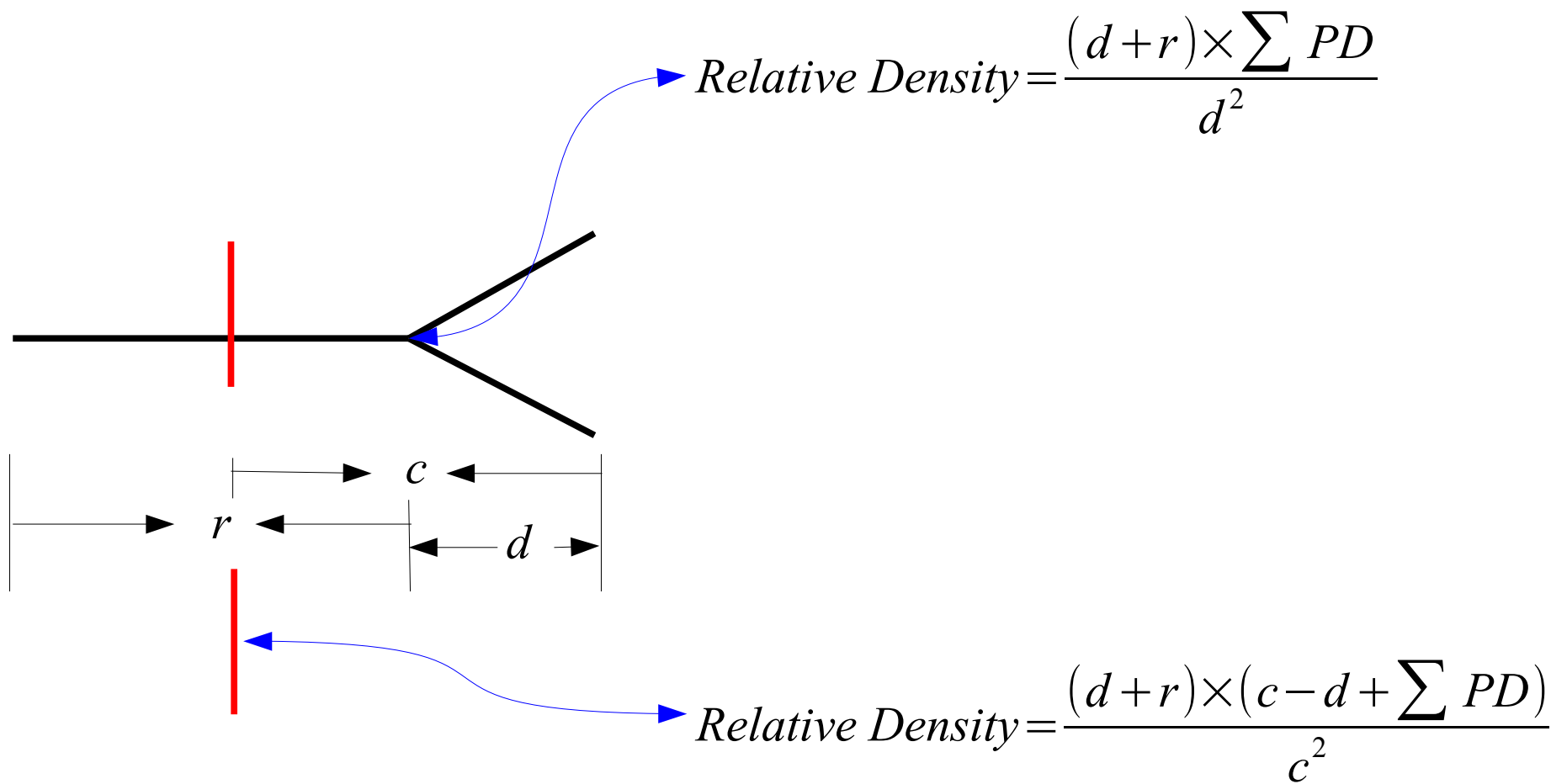
Speed up the OTU script:

Adept data structure for the rooted tree according to the following paper:

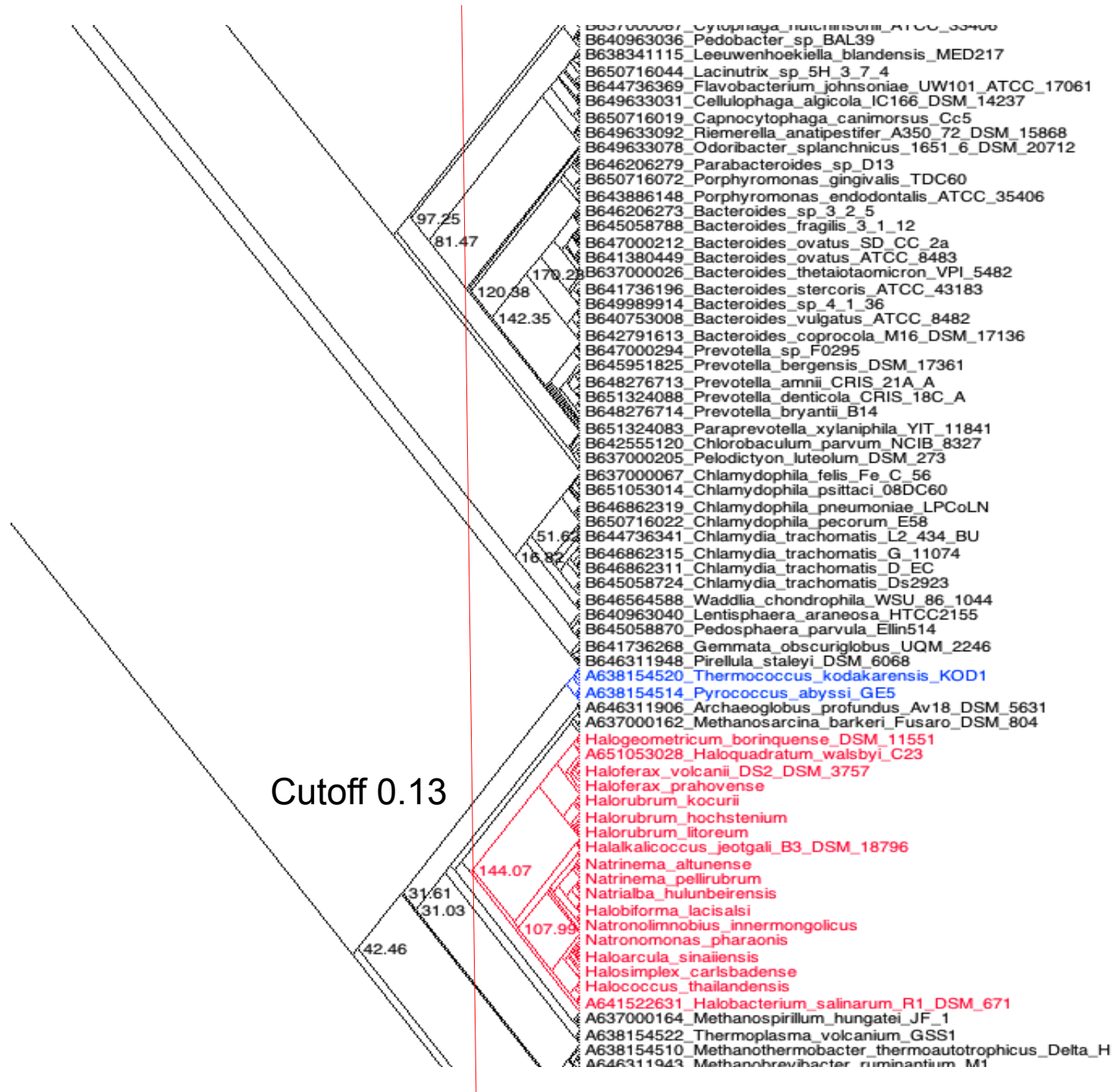
Syst Biol. 2006 Oct;55(5):769-73. Phylogenetic diversity within seconds.
Minh BQ, Klaere S, von Haeseler A.



Node based relative densities are not comparable, densities need to be calculated at the cutoff line instead of the nodes that defines the OTUs

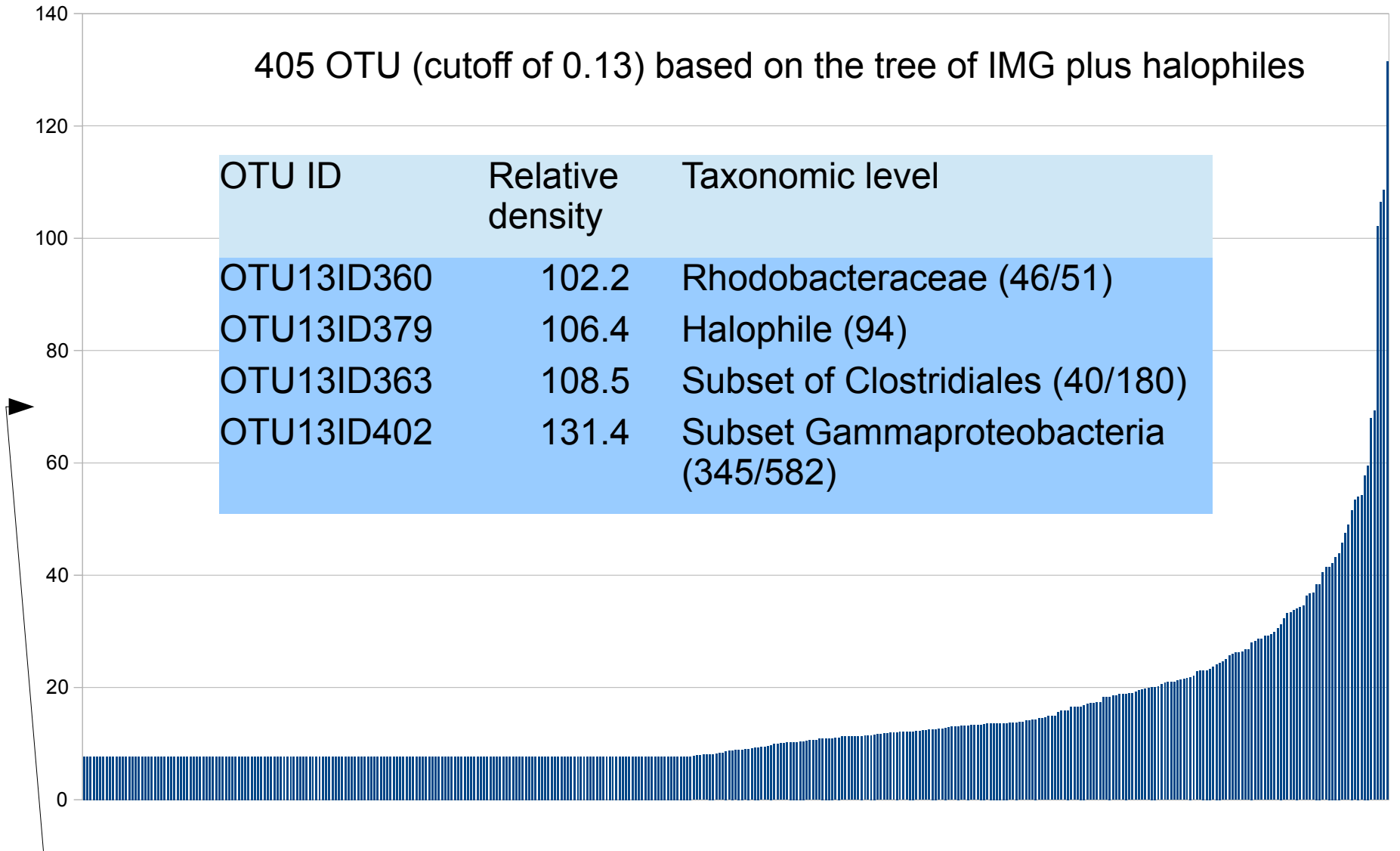


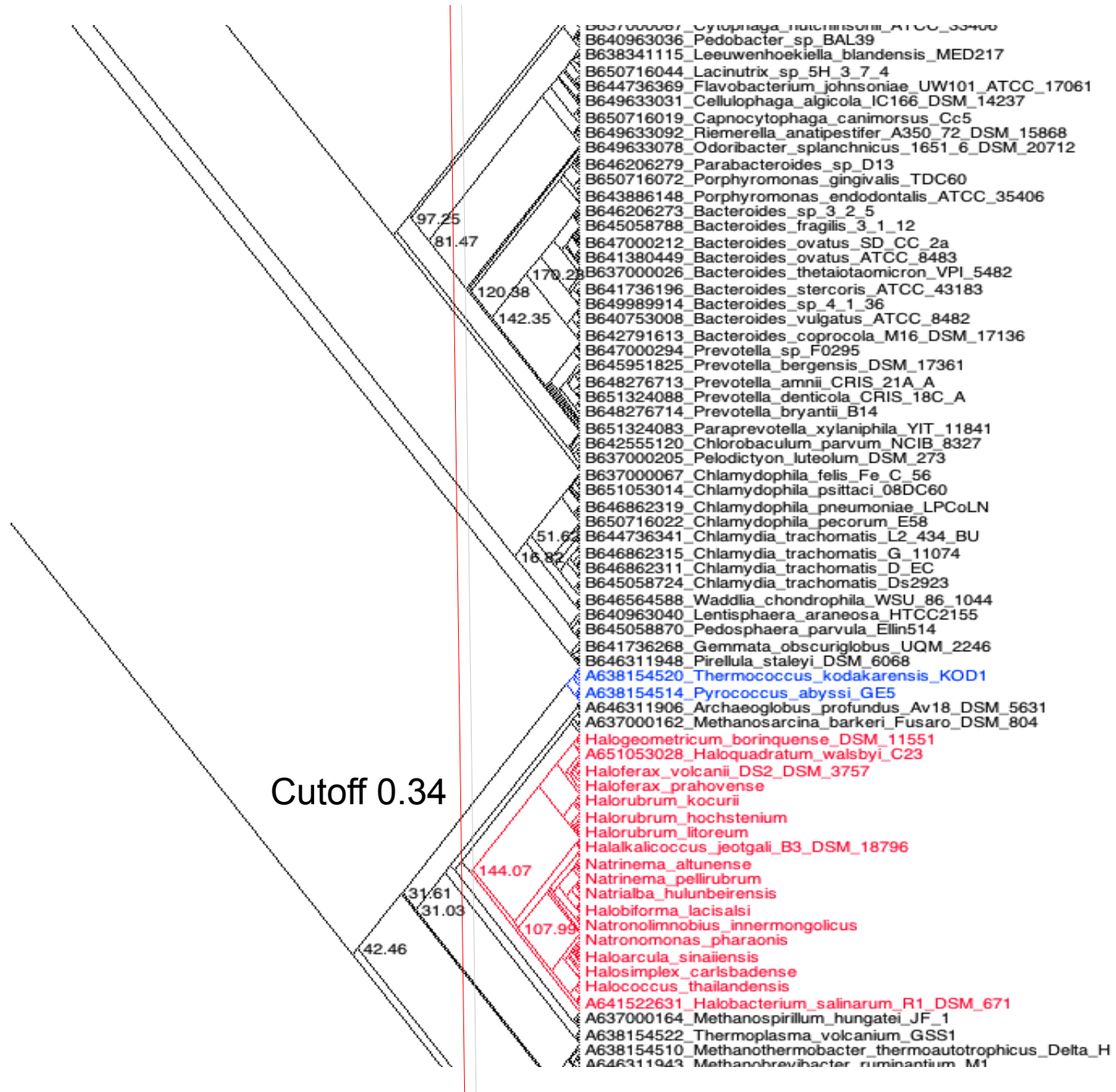
Relative Density at a cutoff line

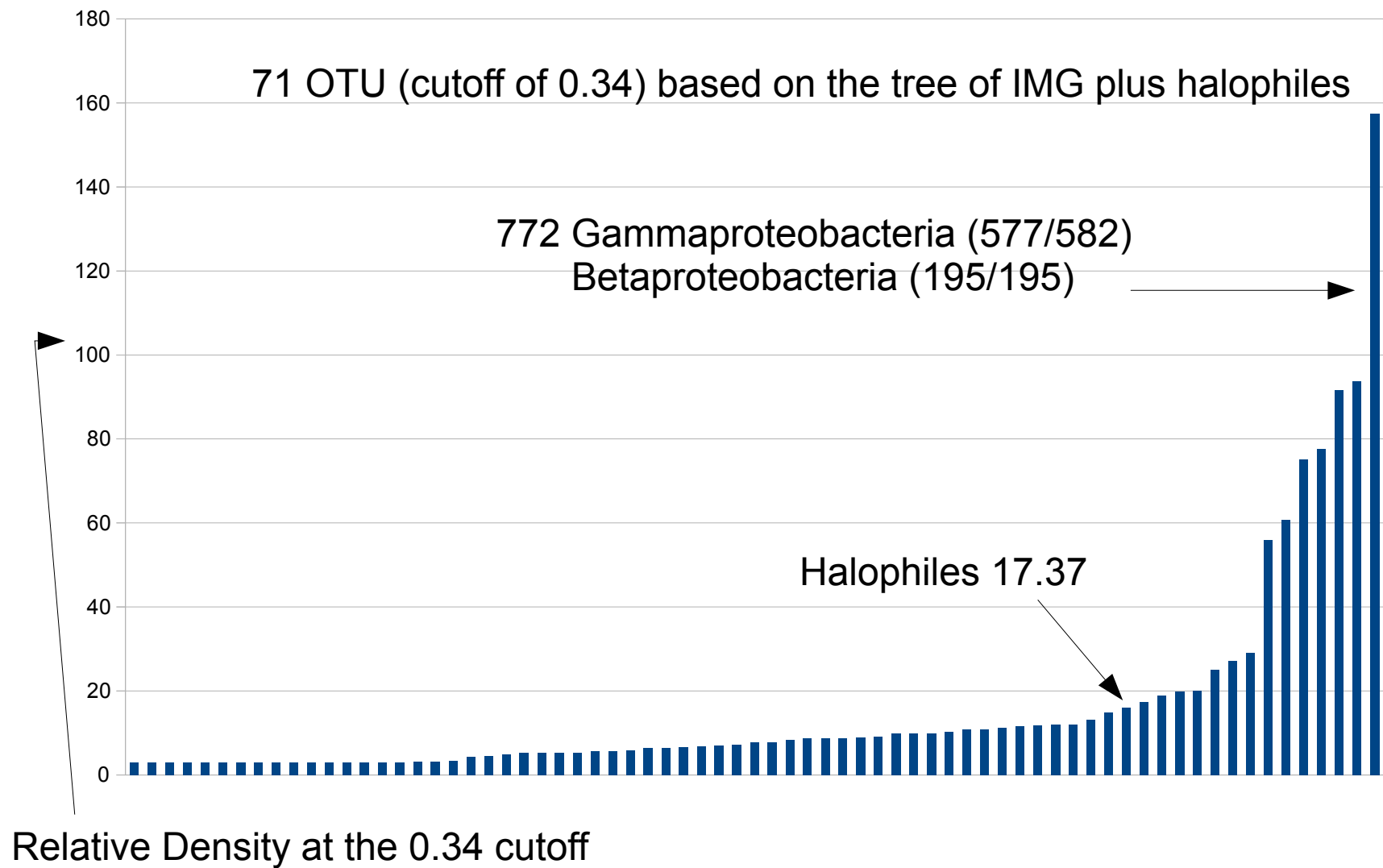


405 OTU (cutoff of 0.13) based on the tree of IMG plus halophiles

OTU ID	Relative density	Taxonomic level
OTU13ID360	102.2	Rhodobacteraceae (46/51)
OTU13ID379	106.4	Halophile (94)
OTU13ID363	108.5	Subset of Clostridiales (40/180)
OTU13ID402	131.4	Subset Gammaproteobacteria (345/582)







Summary of the halophile sample density study:

Based on a fasttree built from all IMG bacteria/archaea with halophile samples, halophile group has a relative density of 106.4 at the 87% relative depth level from the root. It is one of the top 4 best sampled group at this level. From the level of 87% to 66% towards the root, the relative sample density decreases to 17.4 because no more halophile genomes are added.

Test Run of OTU grouping on Merlot: compute-5-10

1. OTU grouping

Input: Greengenes fasttree of 399,817 taxa (mid point rooted)

Cutoff: normalized branch length of 0.10 from the tip

Naive depth calculation:

Start: Mon Jul 9 12:35:21 PDT 2012

End: Mon Jul 9 12:39:49 PDT 2012

Weighted depth calculation:

Tue Jul 10 15:33:01 PDT 2012

Tue Jul 10 15:37:29 PDT 2012

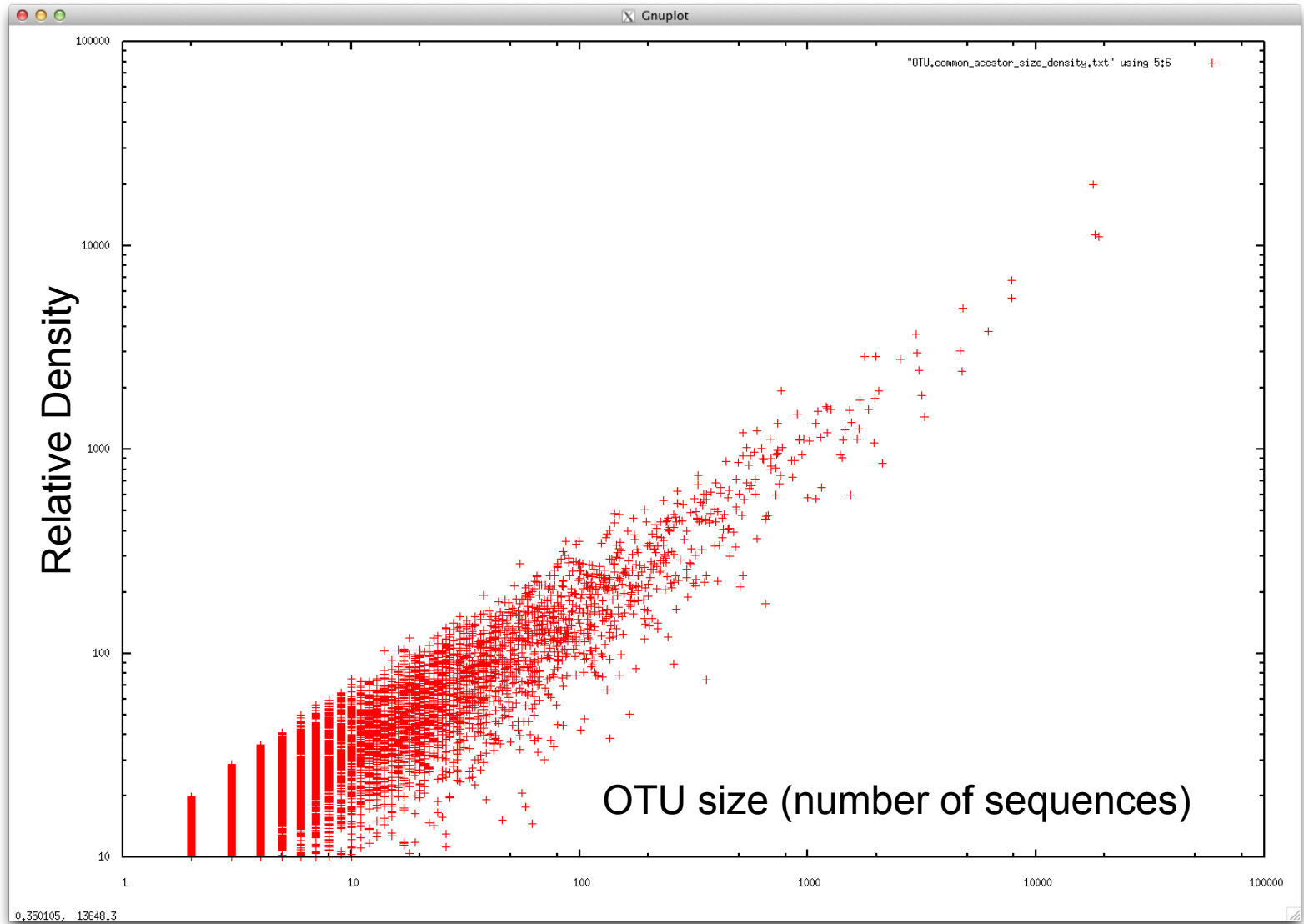
399,817 sequences have been grouped into 52,345 OTUs by naïve average depth method

Using the same cutoff, they are grouped into 69,339 OTUs by weighted average depth method

2. Calculate Density at the cutoff of 0.10 from tip for all the OTU groups

Relative density of OTUs at the 0.10 relative depth level from the tips

Density	OTU count	ssu-rRNA count
10(one seq per OTU)	38781	39453
11-20	7483	21084
21-50	4129	34705
51-100	1112	32180
101-200	477	34626
201-500	254	50525
501-1000	61	37722
1001-2000	33	43390
2001-5000	10	35688
5001-10000	2	15641
10000-20000	3	54803



The most densely sample OTUs (0.10 cutoff)

OTU52330 ssu-rRNA count:17842 Relative Density:19864.4

201 Unclassified:OTU

10601 k__Bacteria:p__Firmicutes:c__Bacilli:o__Bacillales:f__Staphylococcaceae:g__Staphylococcus:s__Staphylococcus:OTU

7037 k__Bacteria:p__Firmicutes:c__Bacilli:o__Bacillales:f__Staphylococcaceae:g__Staphylococcus:Unclassified:OTU

OTU52315 ssu-rRNA count:18153 Relative Density:11353.3

1190 Unclassified:OTU

16376 k__Bacteria:p__Actinobacteria:c__Actinobacteria:o__Actinomycetales:f__Propionibacteriaceae:g__Propionibacterium:s__Propionibacterium:OTU

587 k__Bacteria:p__Actinobacteria:c__Actinobacteria:o__Actinomycetales:f__Propionibacteriaceae:g__Propionibacterium:Unclassified:OTU

OTU52336 ssu-rRNA count:18808 Relative Density:11097.1

9 k__Bacteria:p__Actinobacteria:c__Actinobacteria:o__Actinomycetales:f__Corynebacteriaceae:g__Corynebacterium:s__Brevibacterium:OTU

469 k__Bacteria:p__Actinobacteria:c__Actinobacteria:o__Actinomycetales:f__Corynebacteriaceae:g__Corynebacterium:s__Corynebacterium:OTU

14561 k__Bacteria:p__Actinobacteria:c__Actinobacteria:o__Actinomycetales:f__Corynebacteriaceae:g__Corynebacterium:Unclassified:OTU

3768 Unclassified:OTU

** not all greengenes entries have Hugenholtz taxonomy assignments*

Future direction: An Unifrac alternative

1. OTU grouping reference trees at different cutoff level (**multi-level**)
2. Map metagenomic sequences on the trees
3. Calculate the relative density for the metagenomic subtrees
4. generate density matrix and compare different samples (**intuitive**)