

Exploring the phylogenetic composition of microbiotas

F. Mazel & J. Sanders



ECOSCOPE
routin g the microcosmos

UBC
—
June, 2018



Organisation

Who are we?

Florent Mazel

flo.mazel@gmail.com

Jon Sanders

jonsan@gmail.com

Question

**How can we (better) describe
the composition of microbiomes
using phylogenetic trees?**

Targeted Skills

Reconstruct a microbial phylogenetic tree

Document microbiota compositions using this phylogenetic tree

Statistically test their relationship with metadata

Visualize these compositions

Organisation of the workshop

Theory



Practice



[These Slides]

```
> read.table()
```

**Important point
Try NOT to be shy!**

Do NOT hesitate to ask questions

Important point

Try NOT to be shy!

Do NOT hesitate to ask questions

**Use red/green cards
to inform us about your
coding situation**



INTRODUCTION

Why Exploring the phylogenetic composition of microbiomes?

Why Phylogeny?

What is a phylogeny?

What is a phylogeny?

*Phylogeny:
“a diagrammatic hypothesis
about the history of the evolutionary relationships of a group of organisms”*

What is a phylogeny?

A hypothesis about the relationship between organisms

What is a phylogeny?

A hypothesis about the relationship between organisms

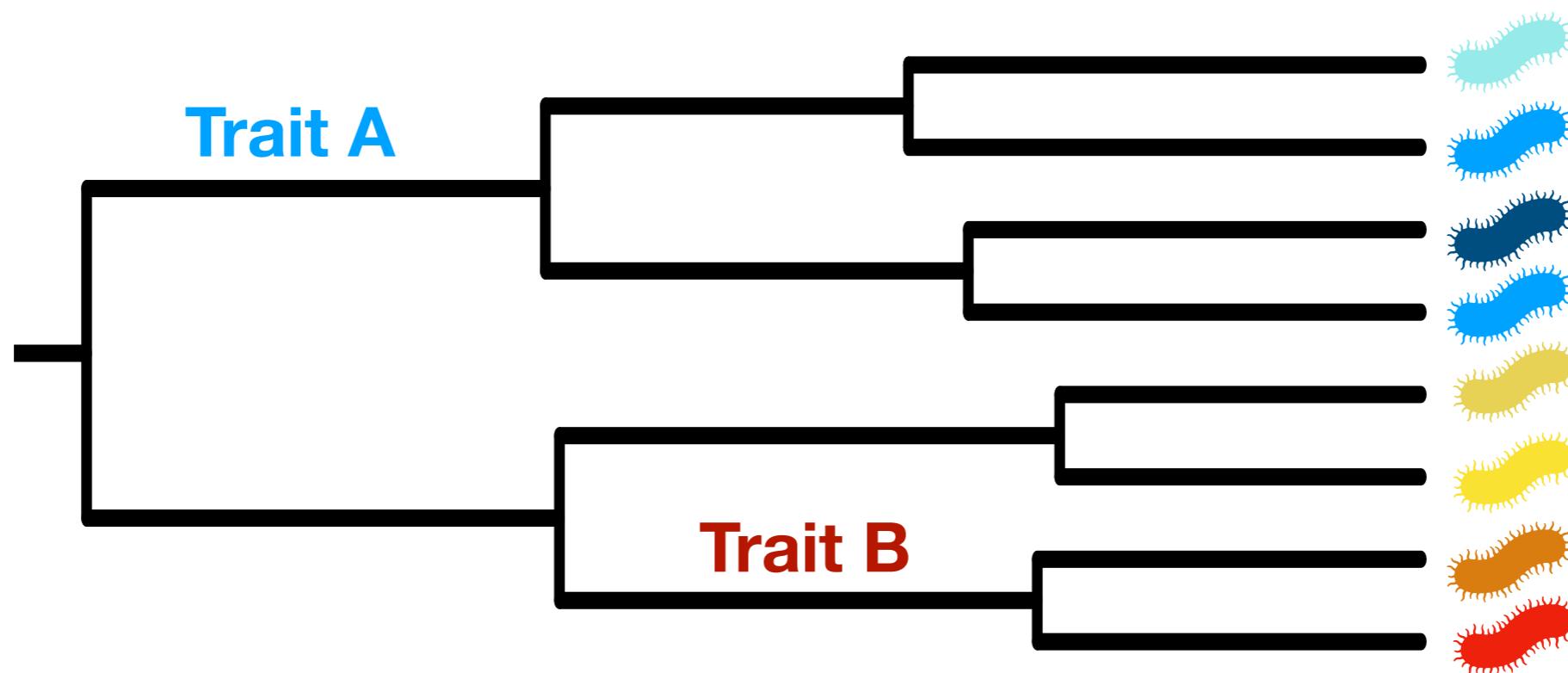
An estimation of something that *actually happened*

What is a phylogeny?

A hypothesis about the relationship between organisms

An estimation of something that *actually happened*
...and hence, a proxy for other things that happened

The phylogeny as a proxy for functions



Phylogeny is a data structure

Taxonomy is also a data structure: systematic categorization



...but most microbes **don't have names**

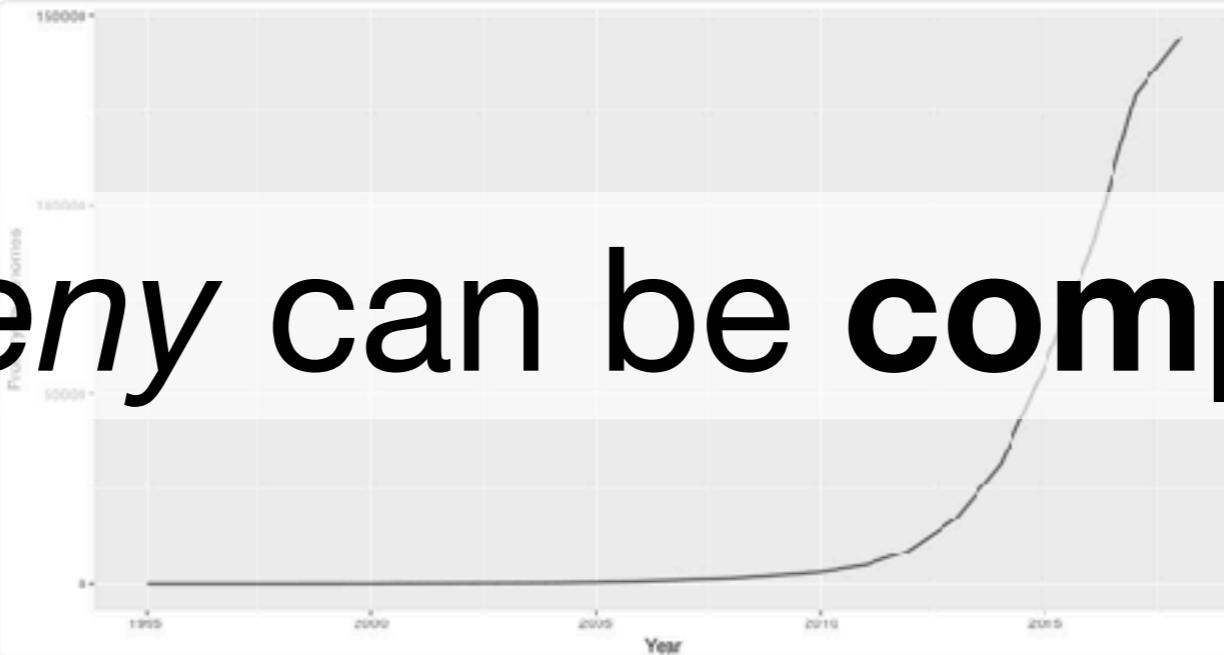
Phylogeny is a data structure

Applications Messages  bacterial genomes

 **Bethany Jose**
@biobeth

Follow 

was looking for a "number-of-bacterial-genomes-over-time" pic for thesis material but most were from ~2014-15, so i made a new one from genbank prokaryotes.txt. number of genomes has nearly doubled during my phd so far. 

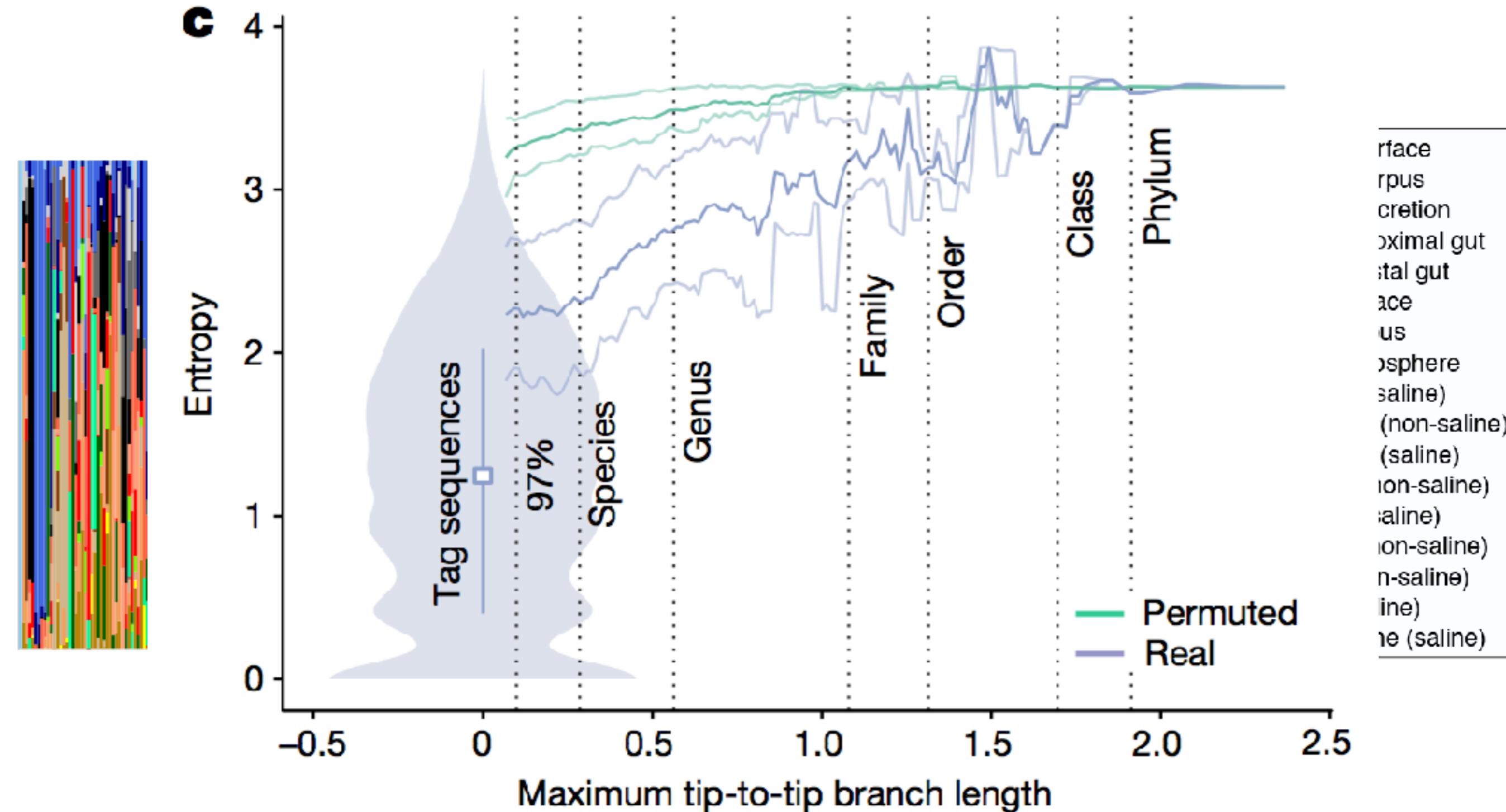


11:14 PM - 13 Jun 2018

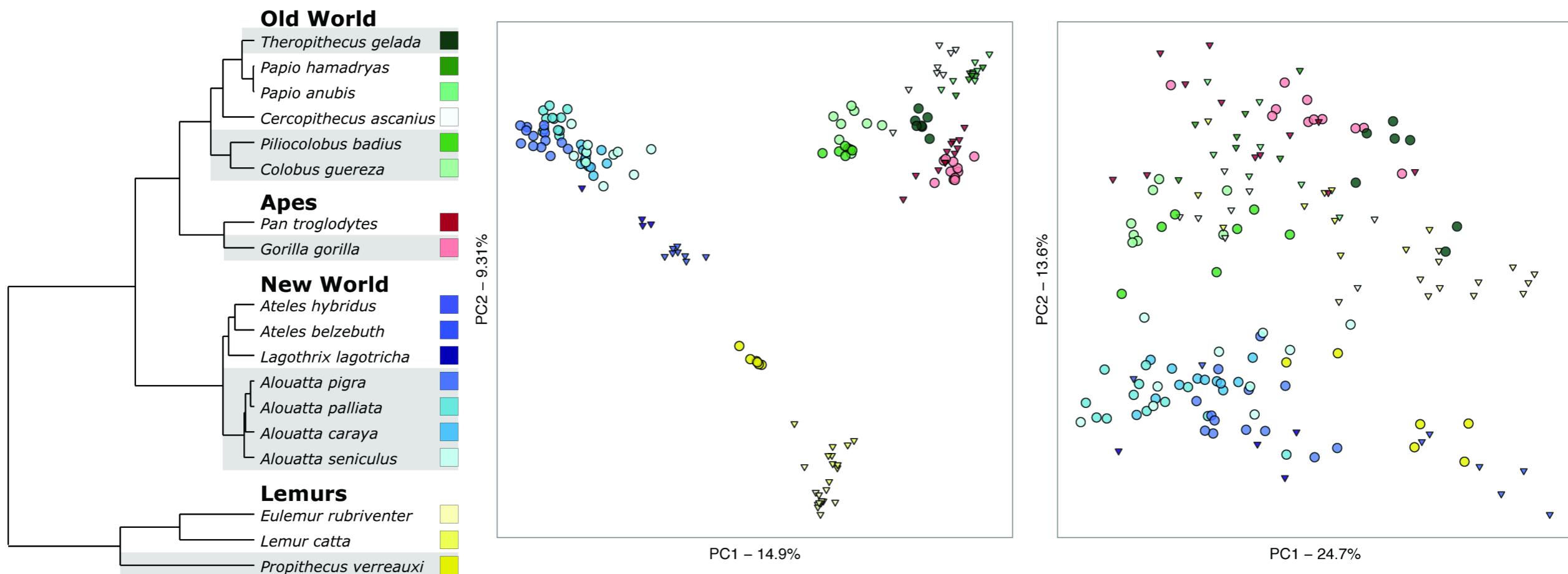
Phylogeny can be computed

Phylogeny correlates well with nature

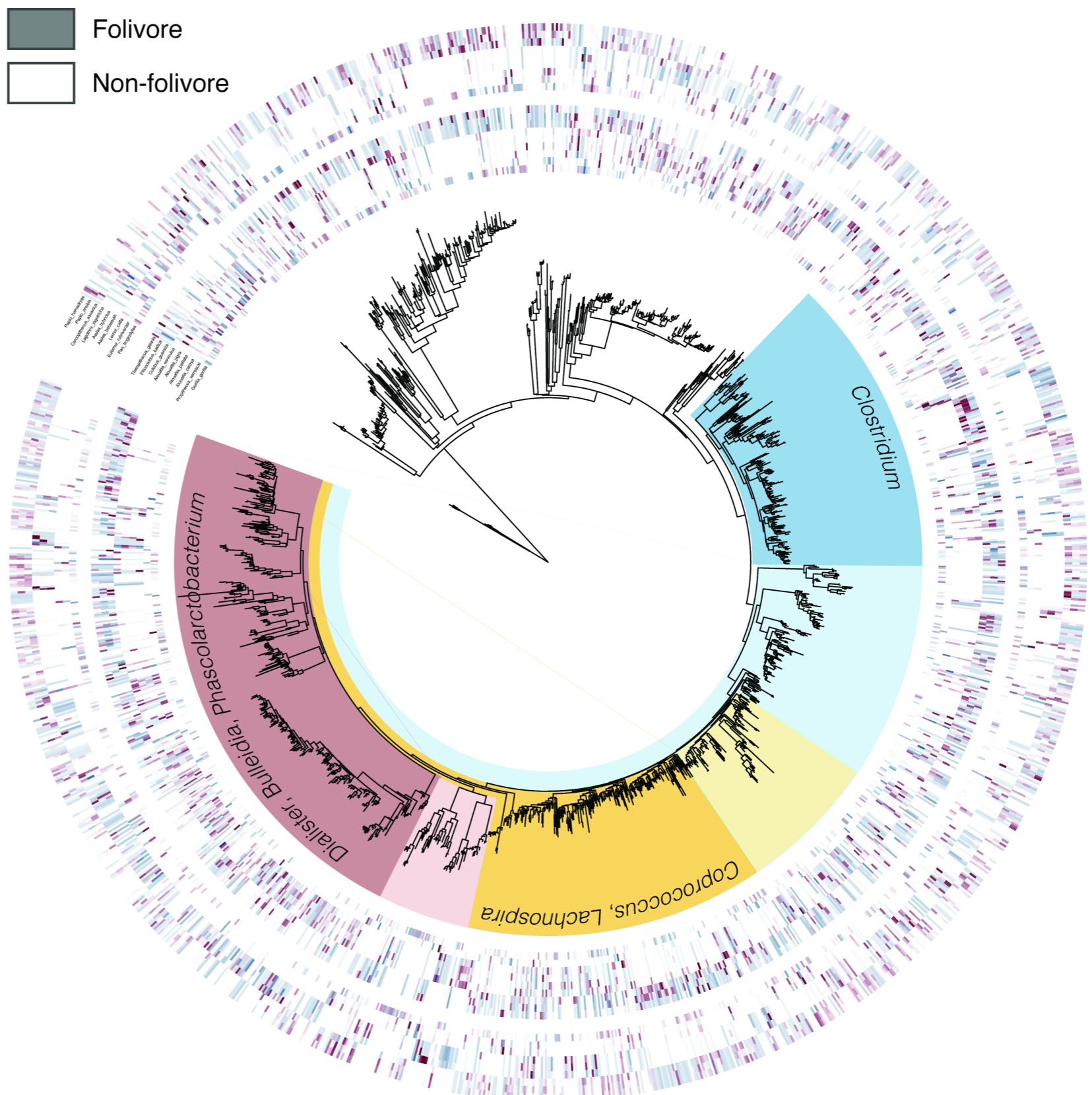
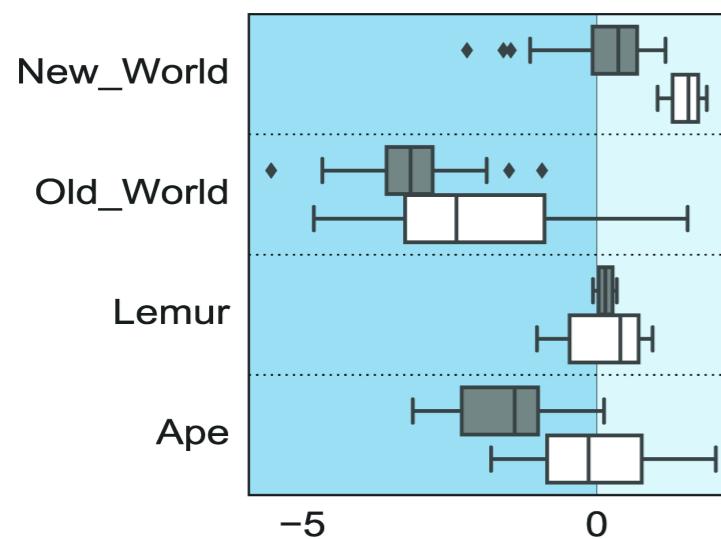
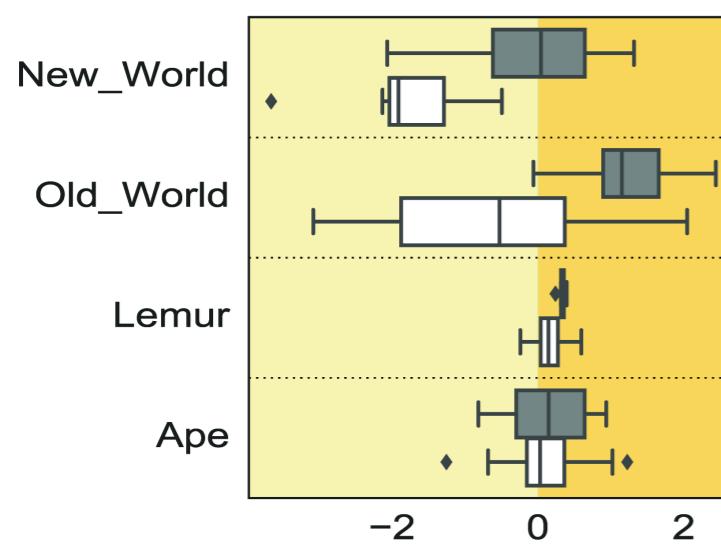
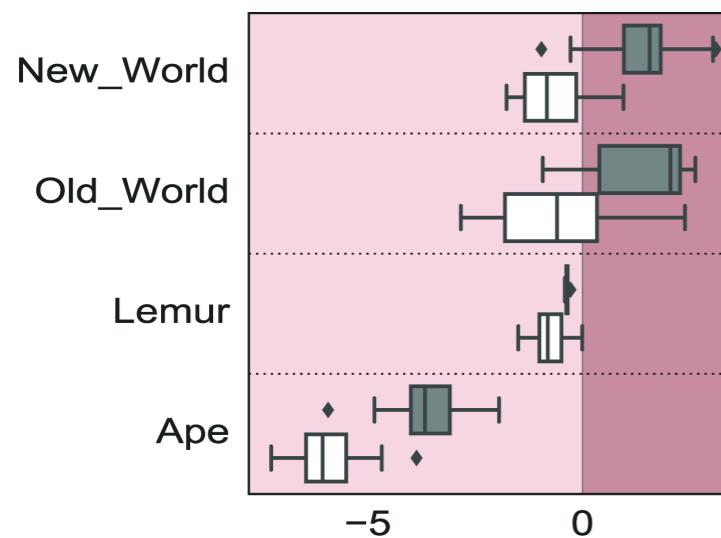
Example: which environments are microbes found in?



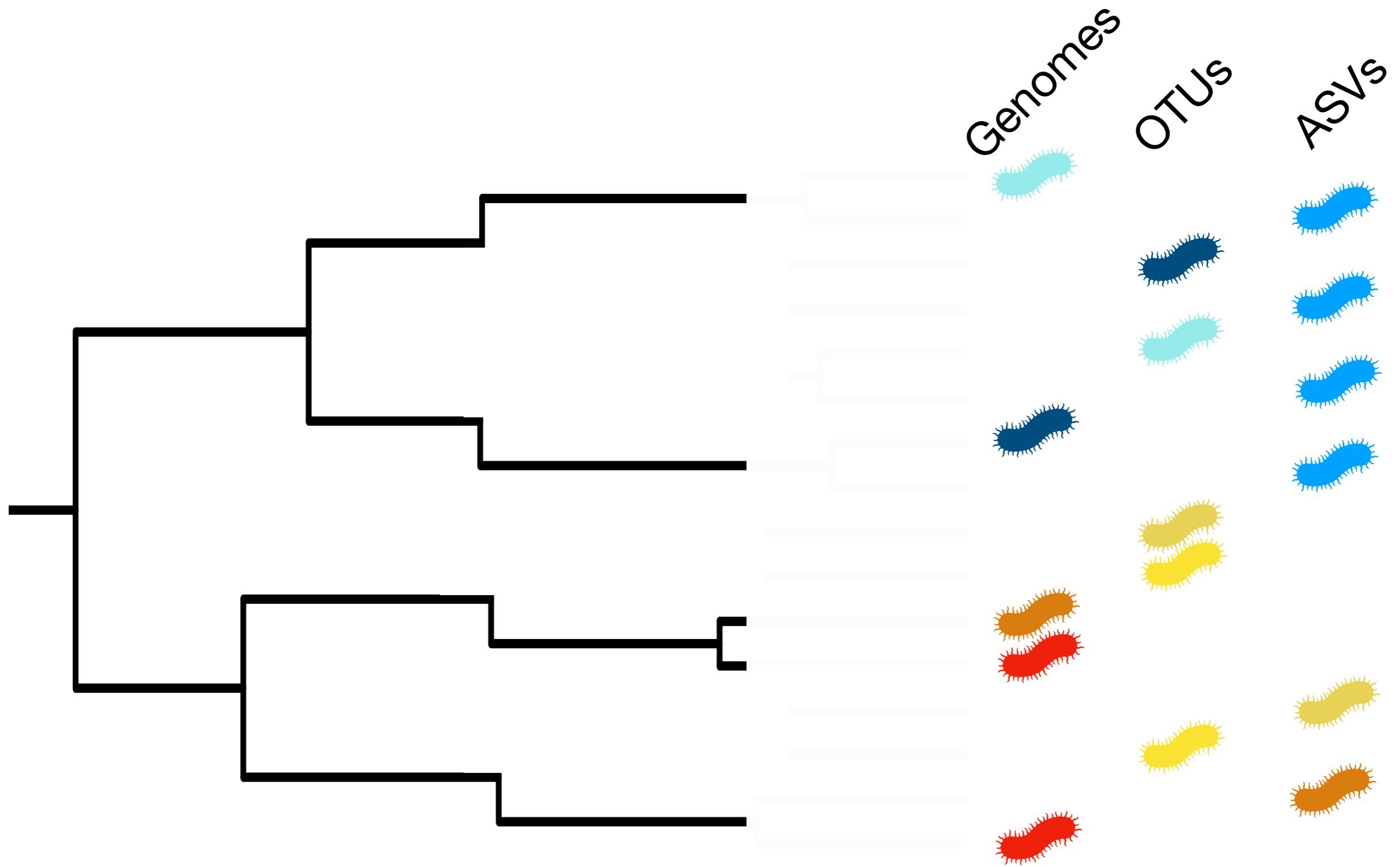
Phylogeny allows scale-independent analysis



Phylogeny allows scale-independent analysis



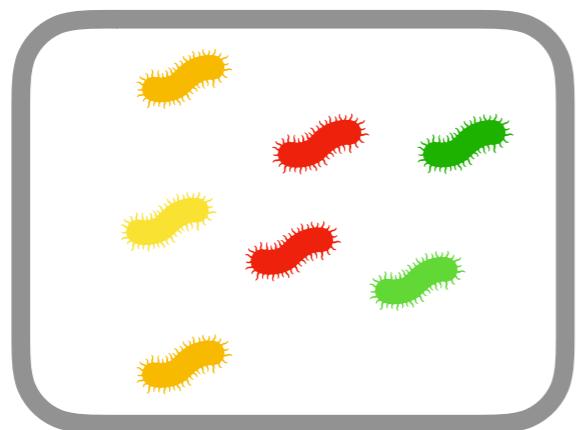
Phylogeny allows integration of data types



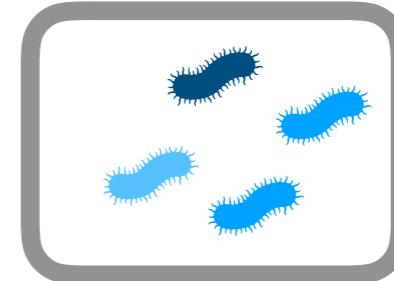
The composition of microbiomes



Microbiome 1

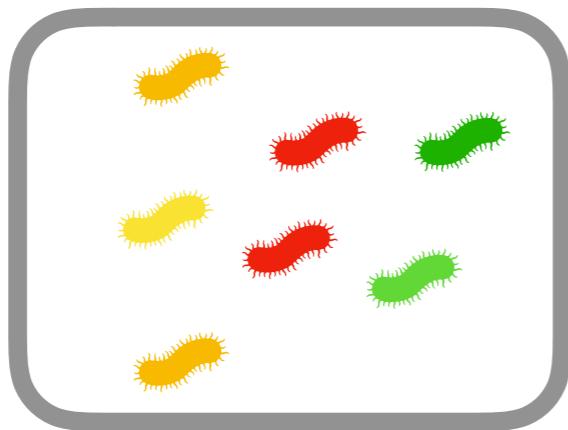


Microbiome 2

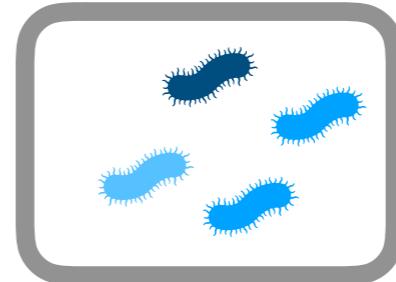


The composition of microbiomes

Microbiome 1

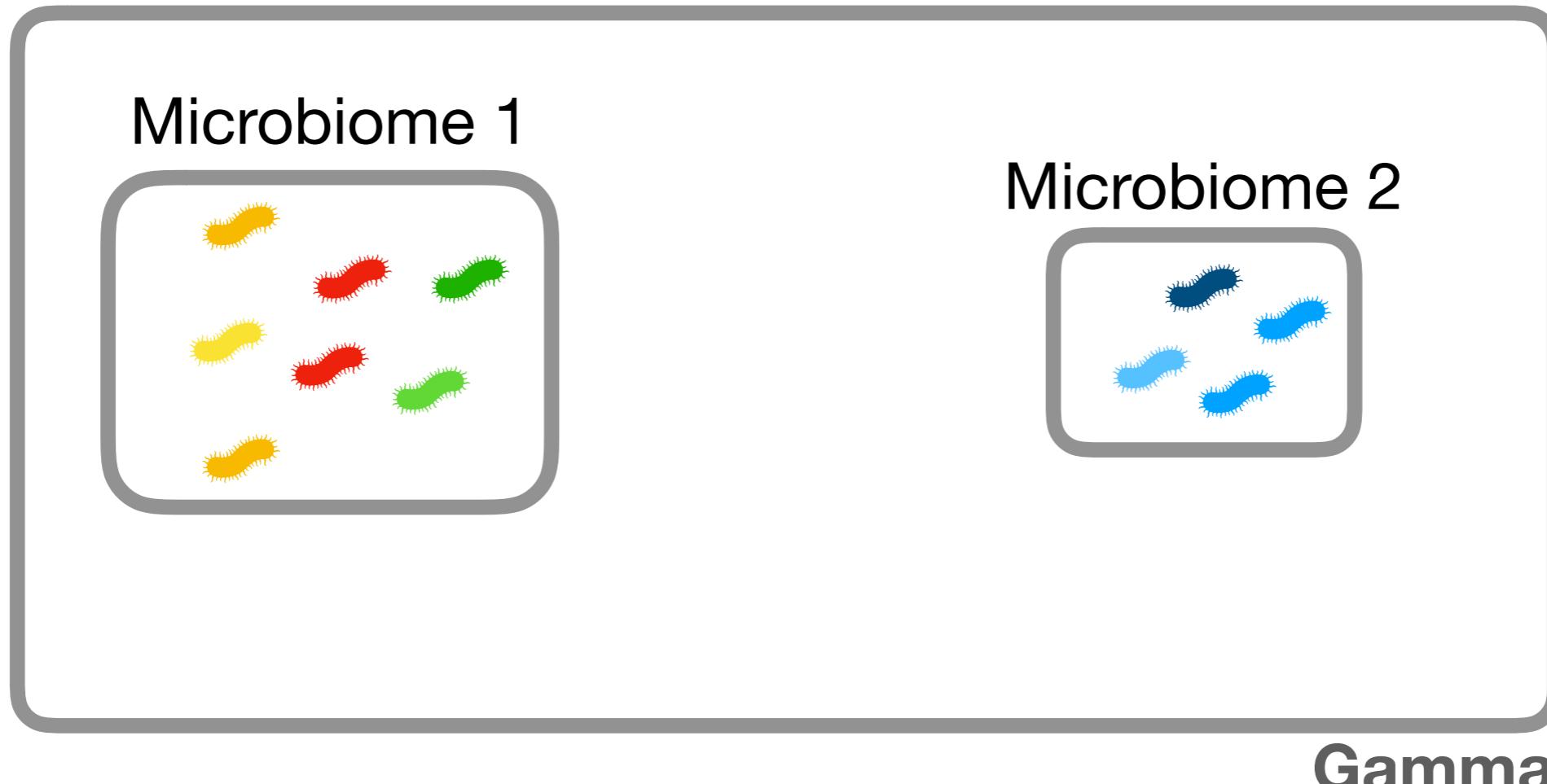


Microbiome 2



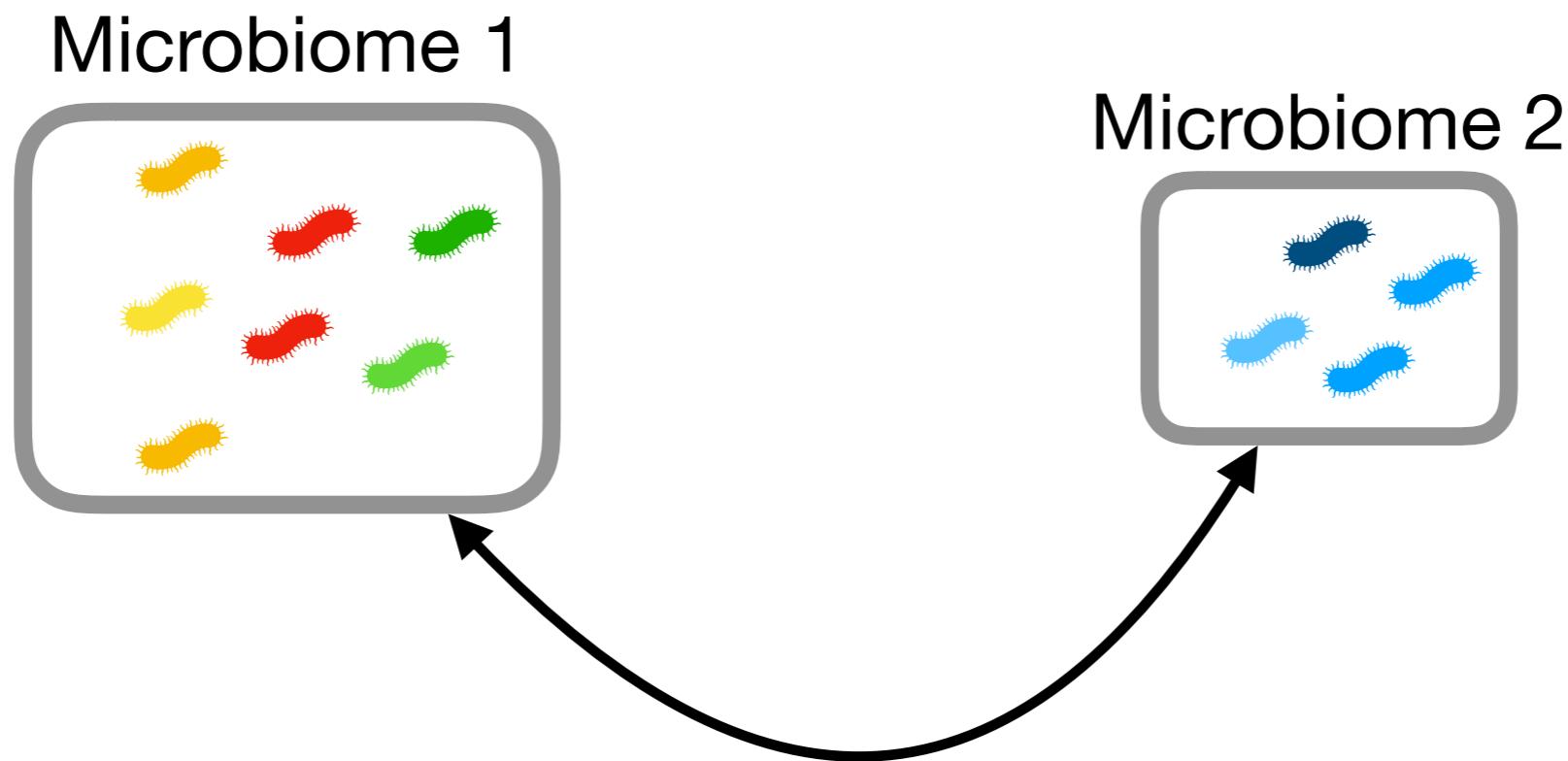
Alpha Diversity = “quantity” of “species” of **ONE** microbiome

The composition of microbiomes



Gamma Diversity = “quantity” of “species” of **ALL** microbiome

The composition of microbiomes



Beta Diversity = “quality” of “species” in the microbiomes

The composition of microbiomes

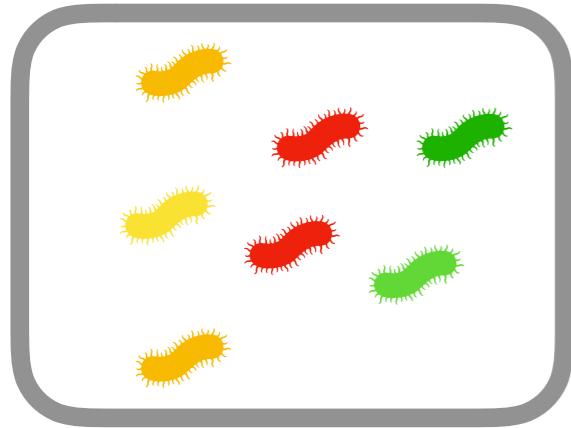
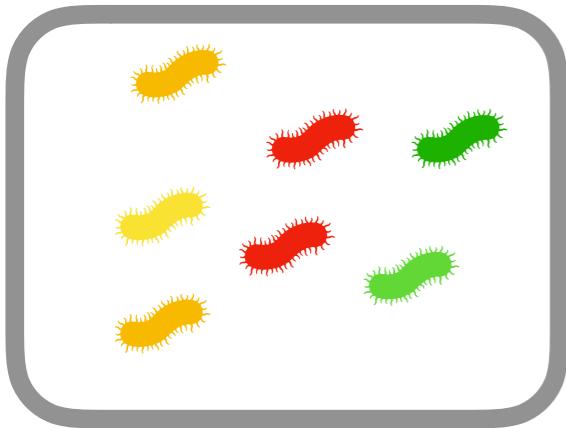
$$\text{Beta Diversity} = \frac{\text{Alpha}}{\text{Gamma}}$$

The composition of microbiomes

$$\text{Beta Diversity} = \frac{\text{Alpha}}{\text{Gamma}}$$

Example

Microbiome 2



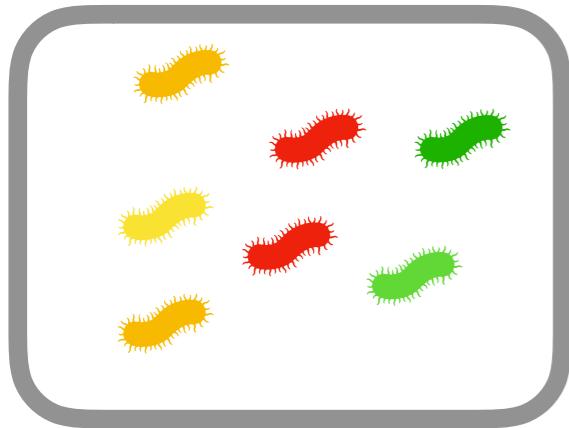
Microbiome 1

The composition of microbiomes

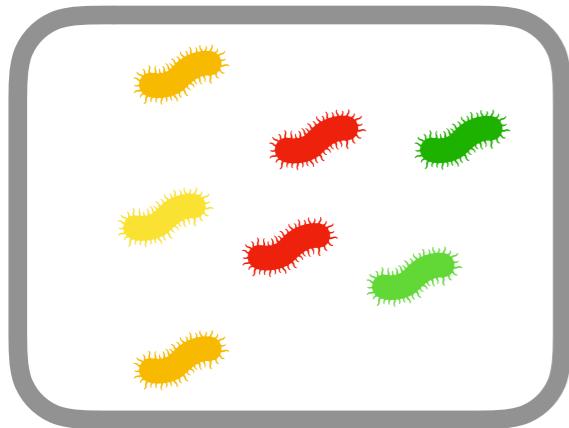
$$\text{Beta Diversity} = \frac{\text{Alpha}}{\text{Gamma}}$$

Example

Microbiome 2



Alpha = 5



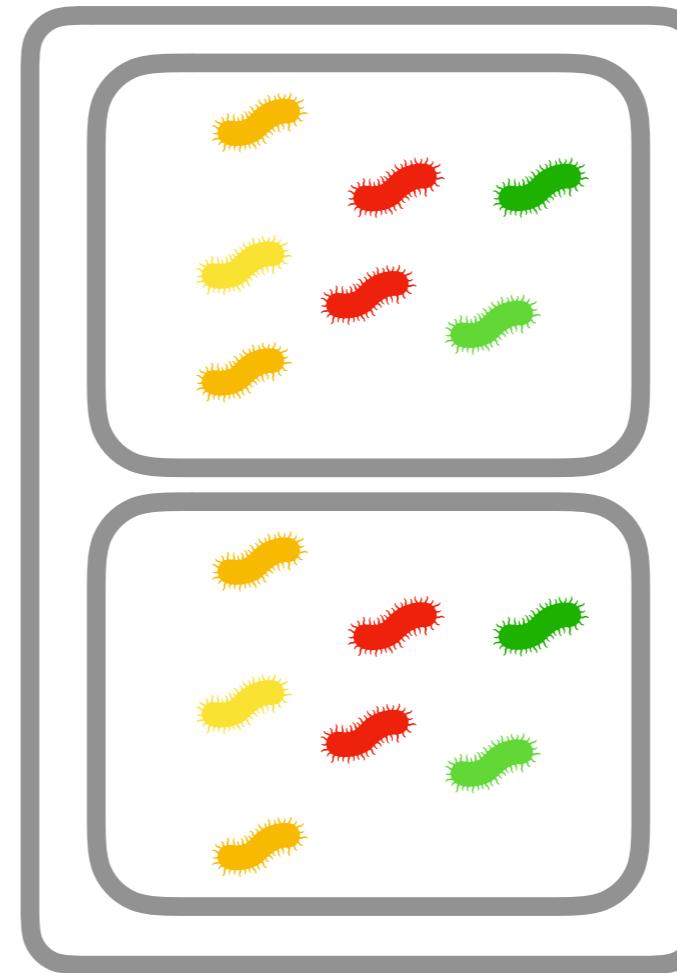
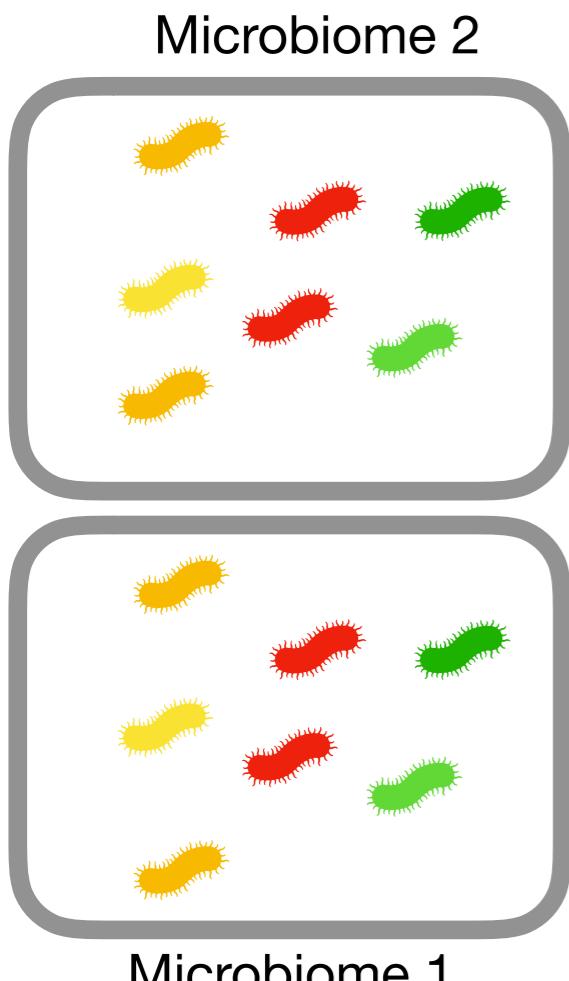
Alpha = 5

Microbiome 1

The composition of microbiomes

$$\text{Beta Diversity} = \frac{\text{Alpha}}{\text{Gamma}}$$

Example

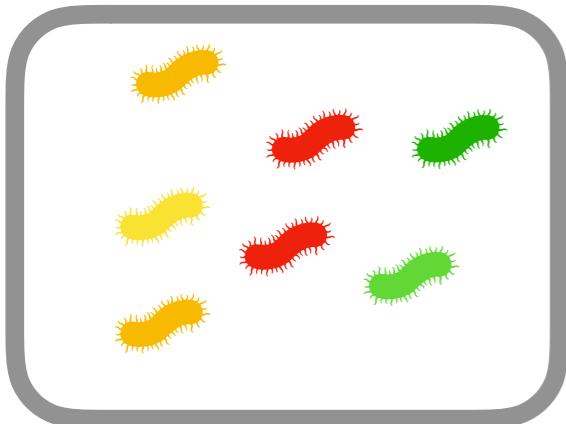


The composition of microbiomes

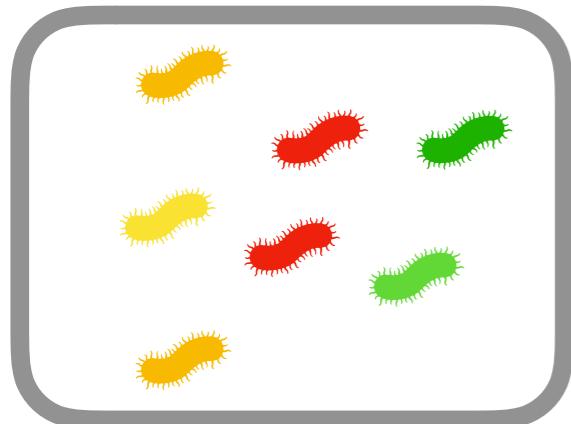
$$\text{Beta Diversity} = \frac{\text{Alpha}}{\text{Gamma}}$$

Example

Microbiome 2



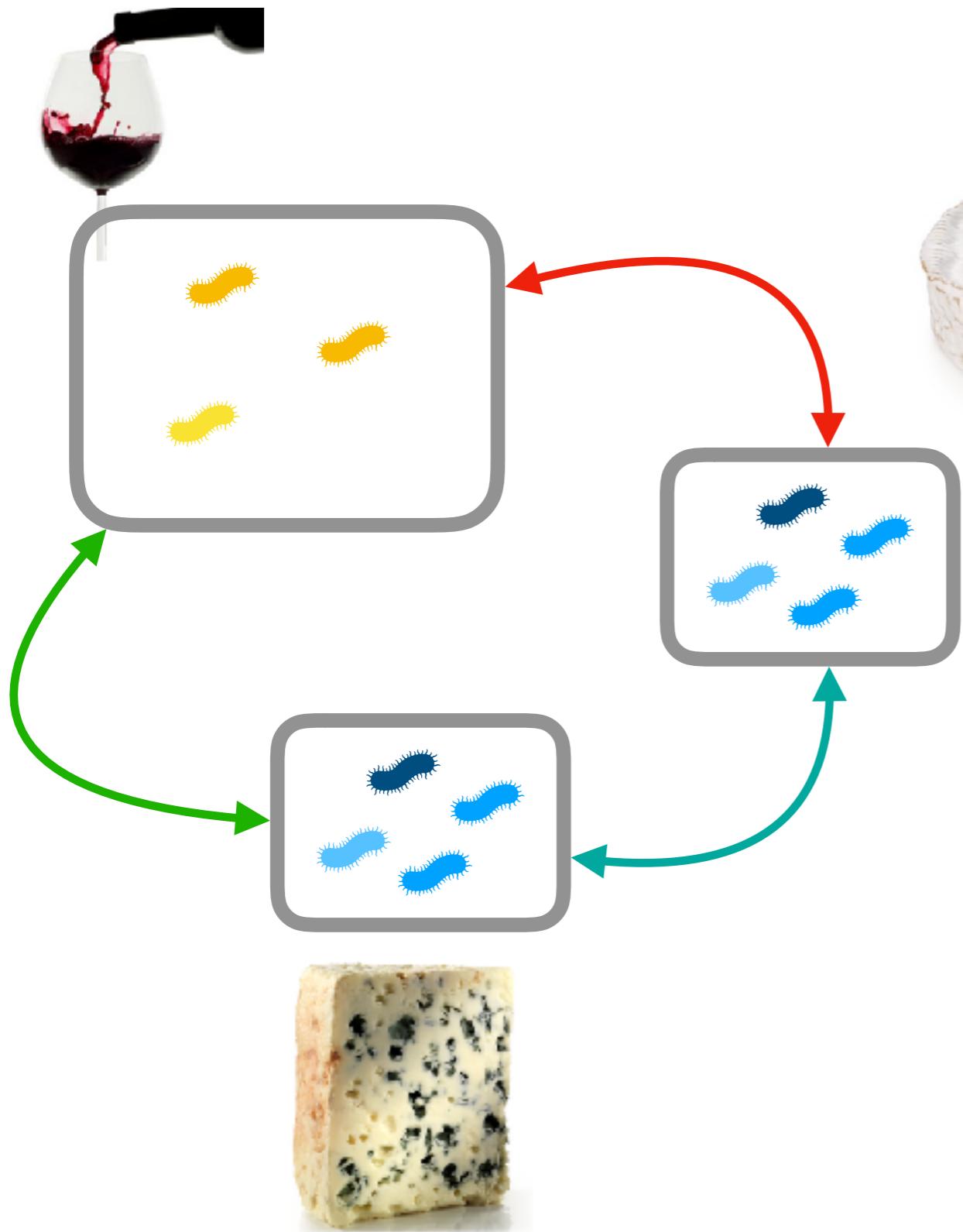
Alpha = Gamma



Beta is minimal

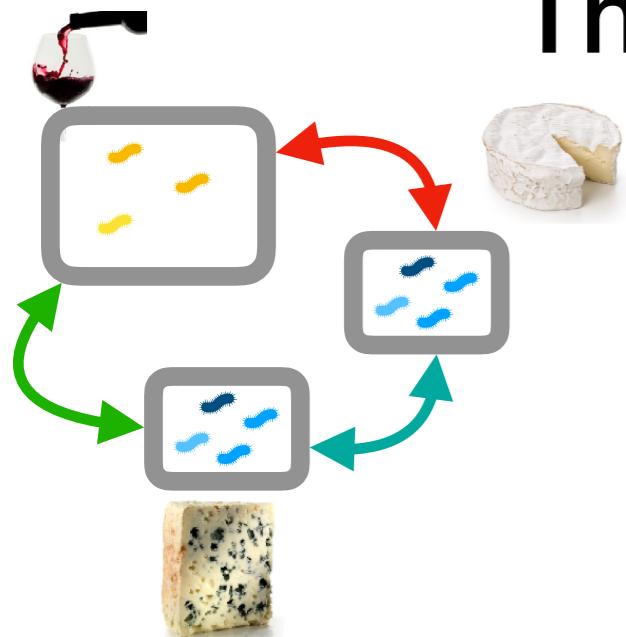
Microbiome 1

The composition of microbiomes

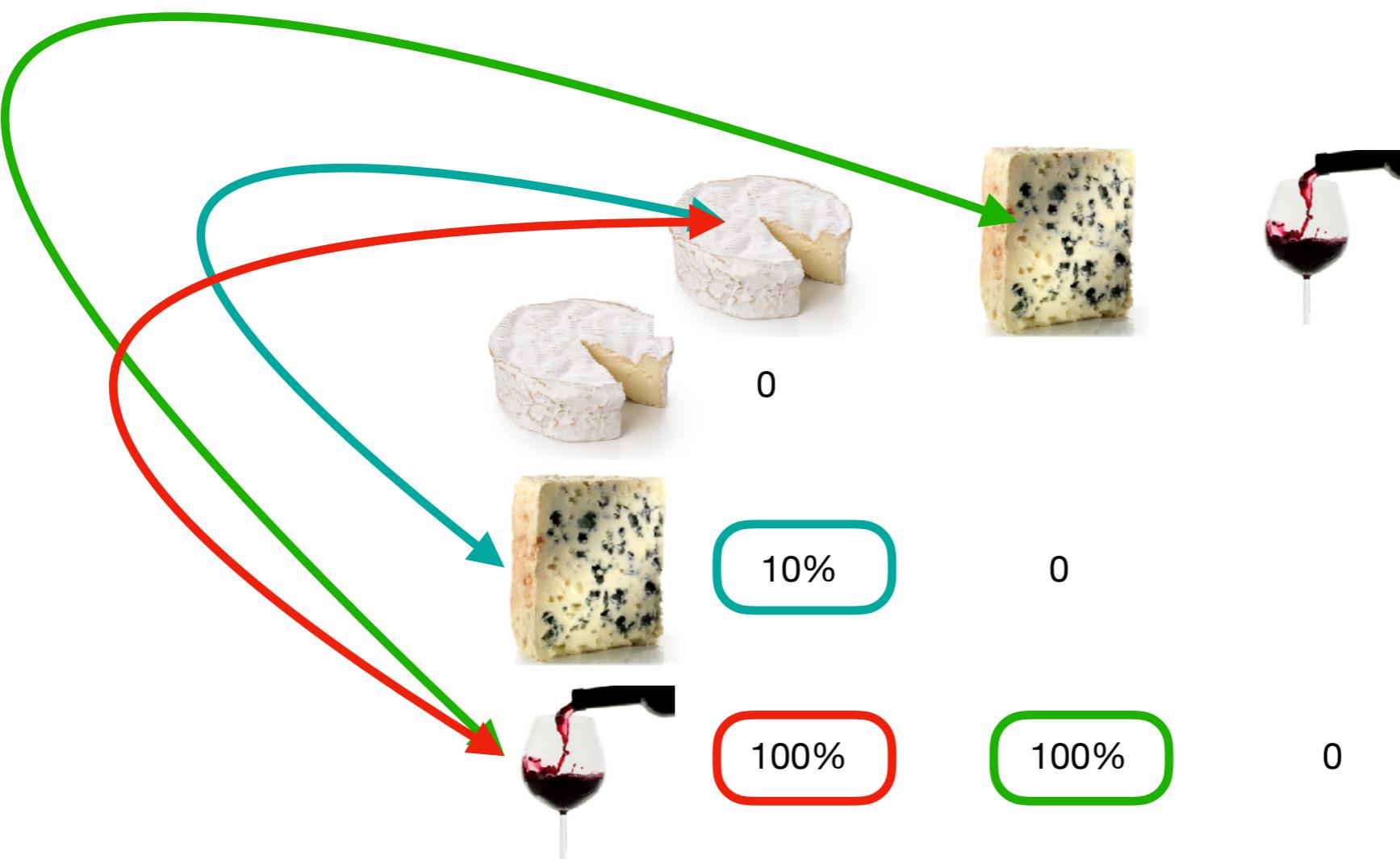


= pairwise measures
of distance or dissimilarities
between microbiota

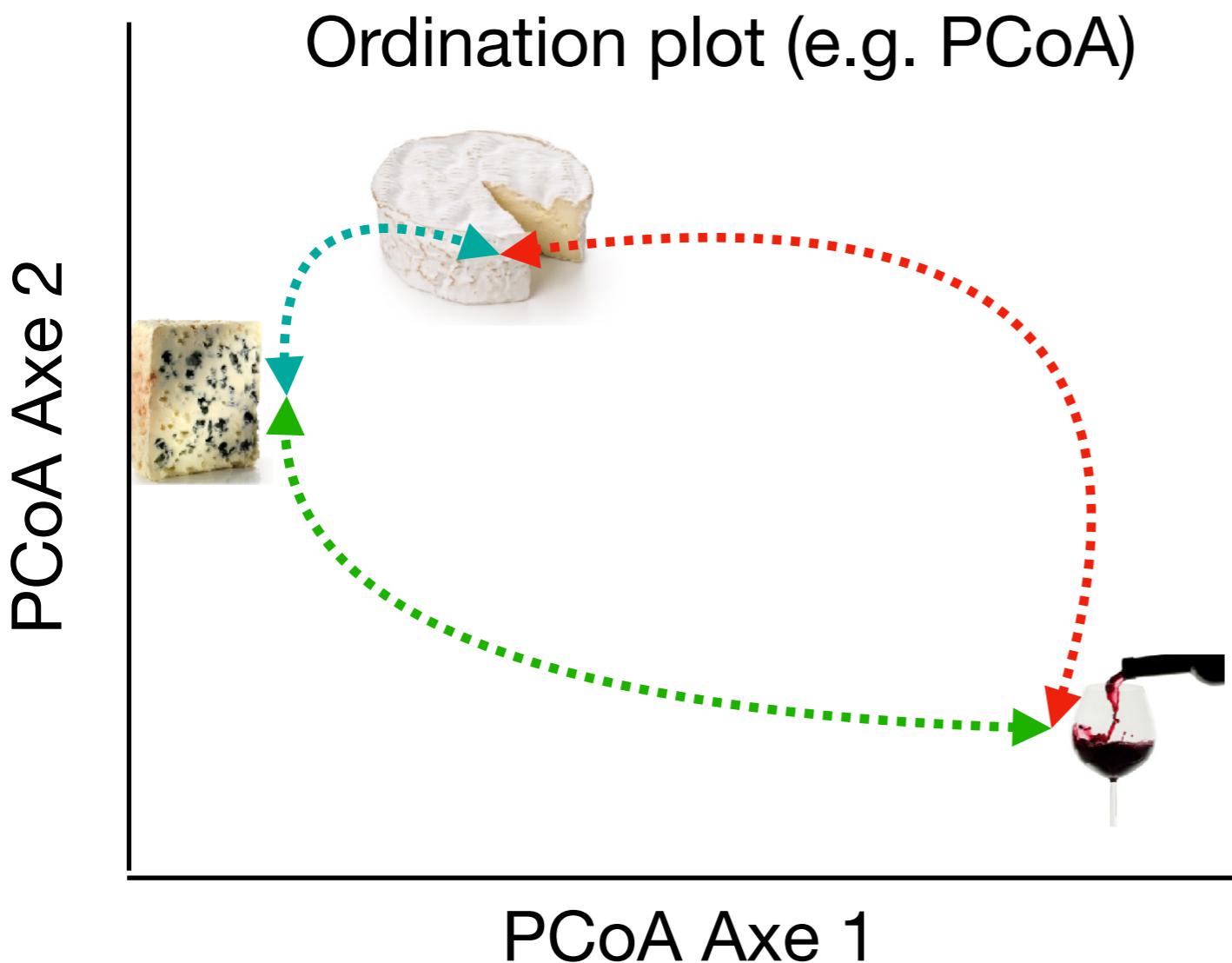
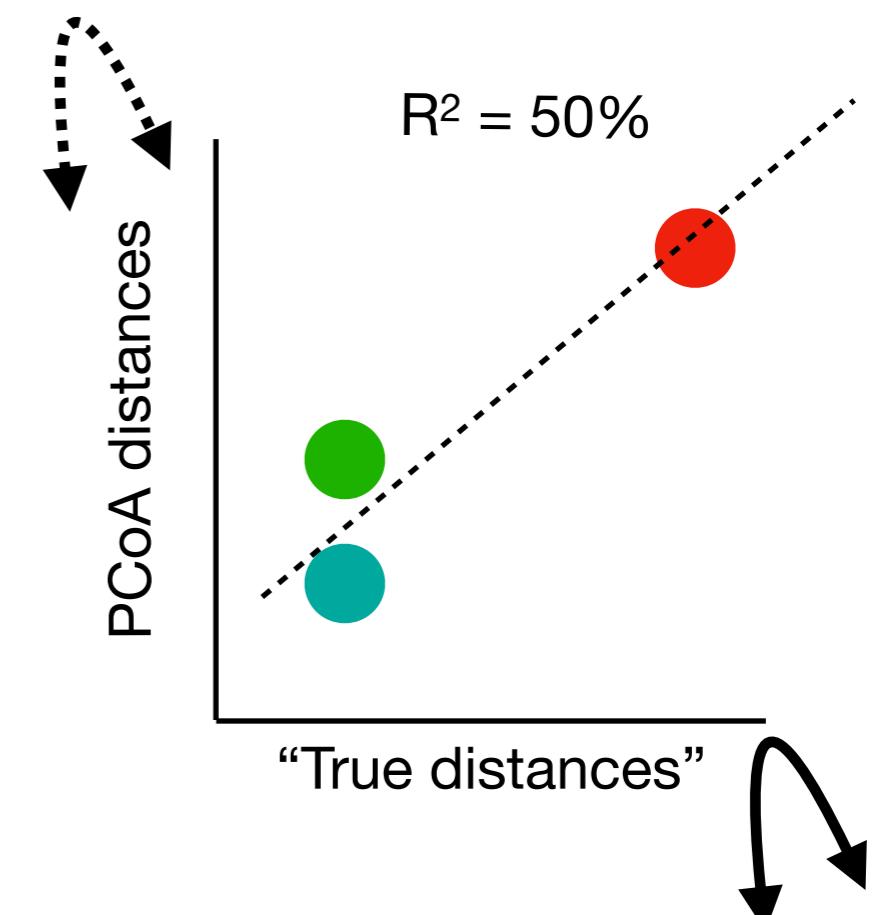
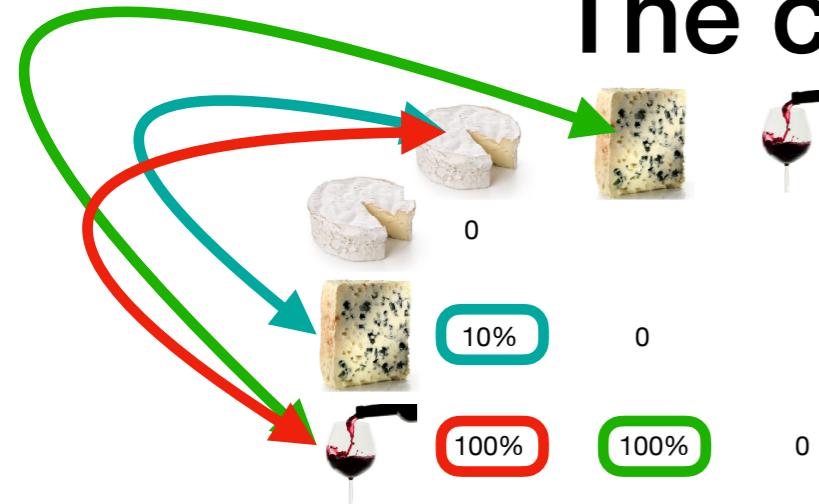
The composition of microbiomes



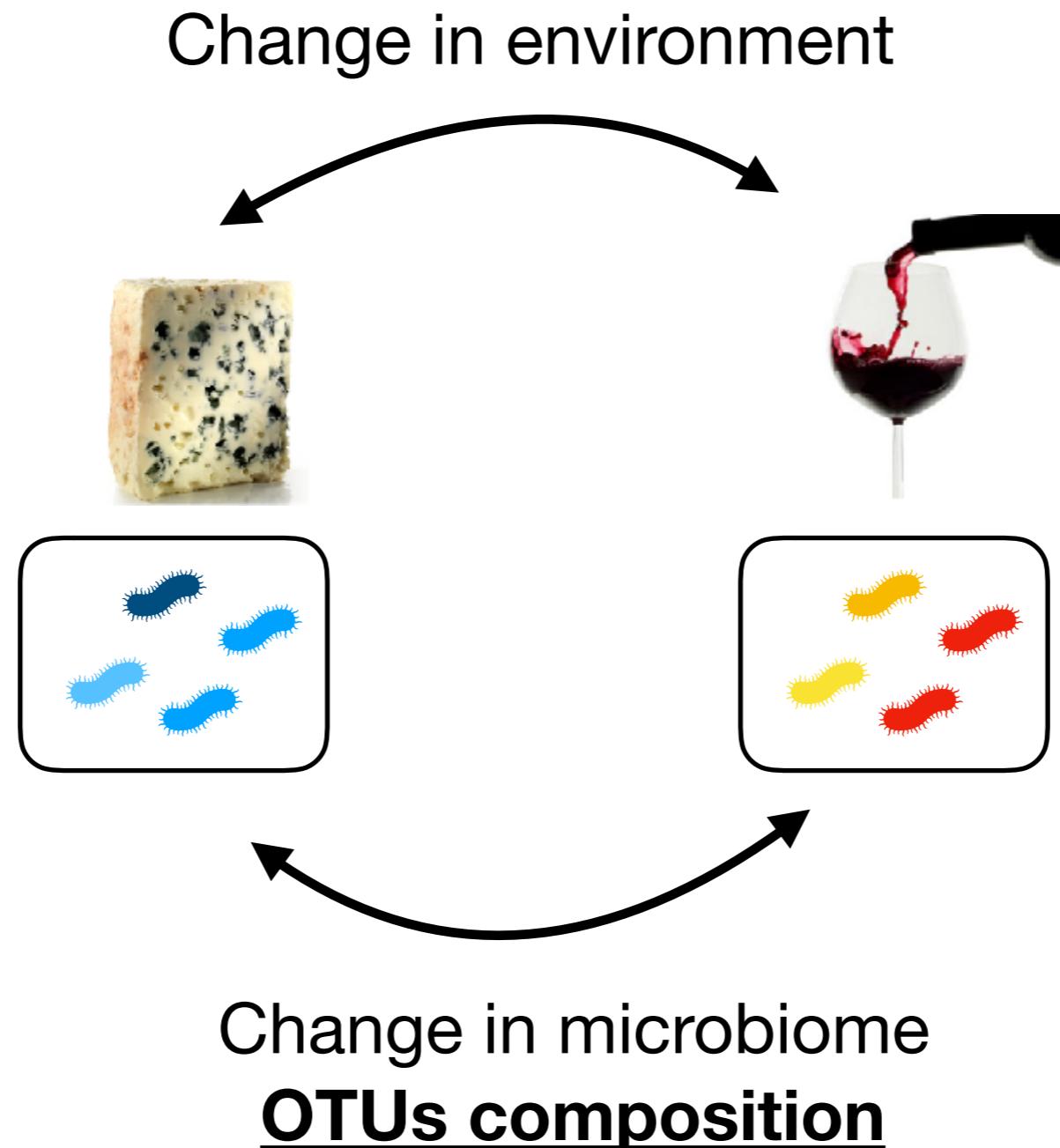
Distance matrix



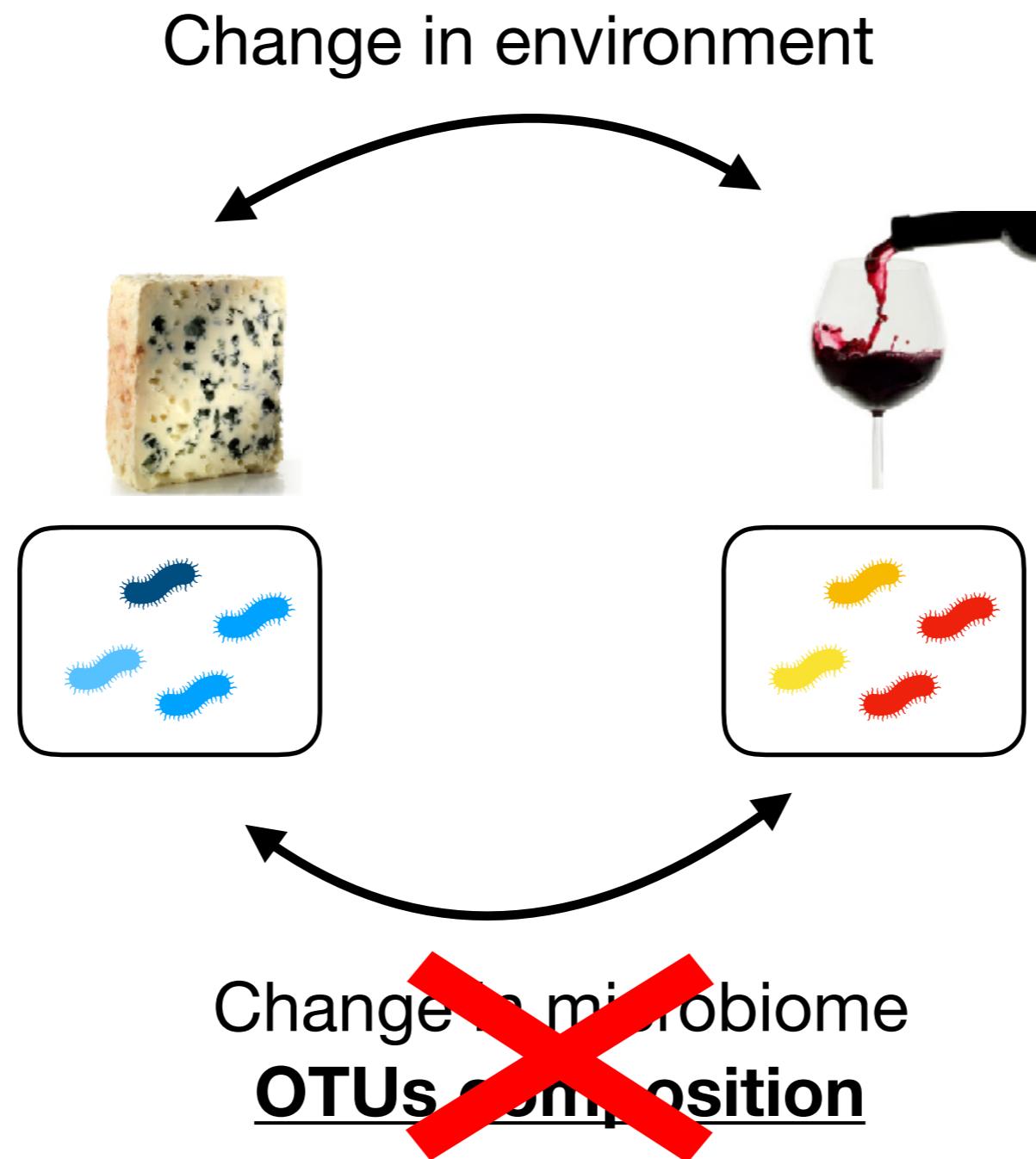
The composition of microbiomes



The composition of microbiomes

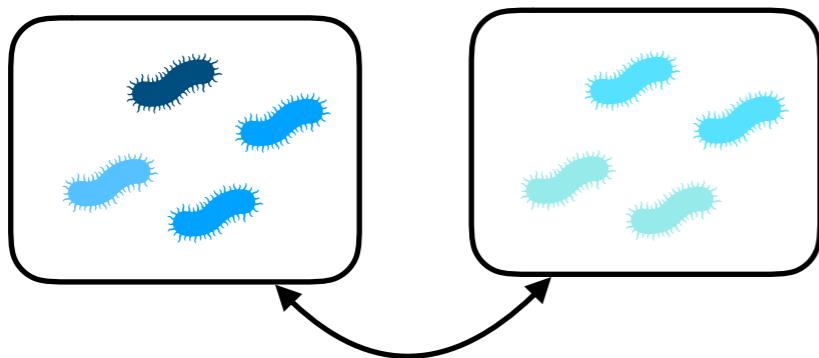


The composition of microbiomes

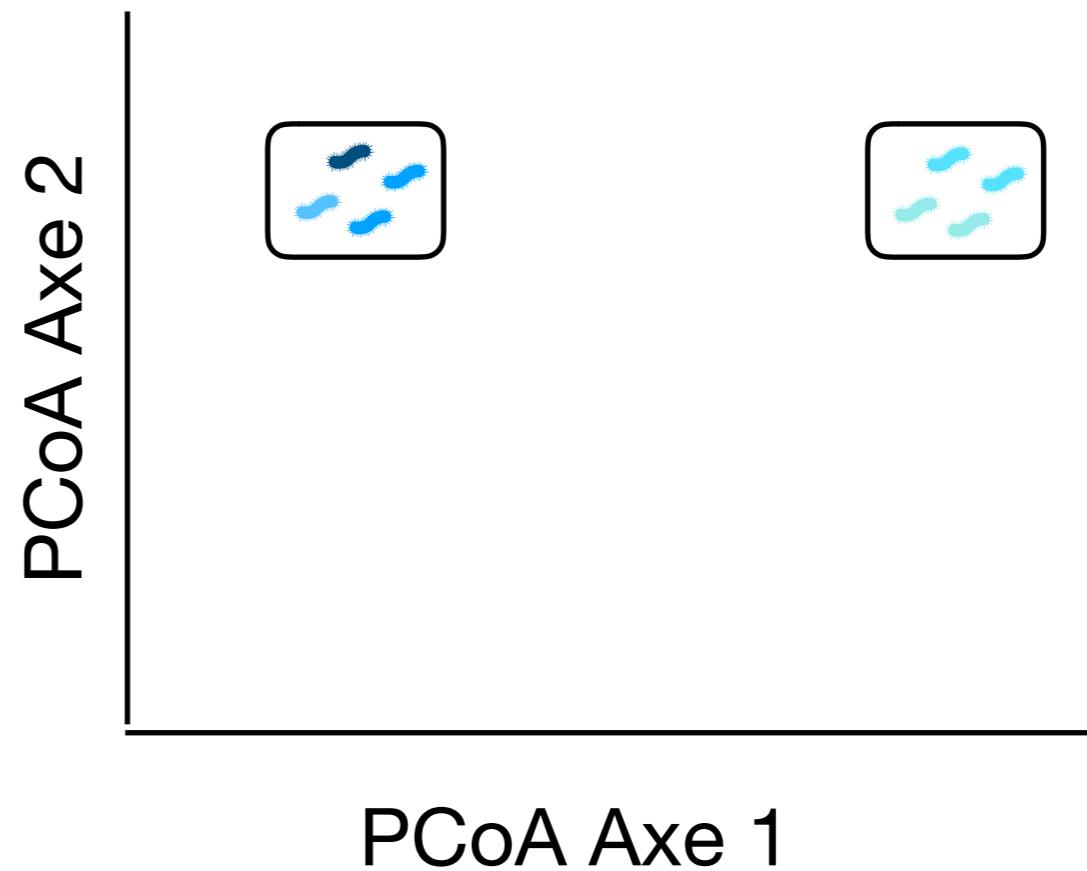


Change in microbiome **functional composition**

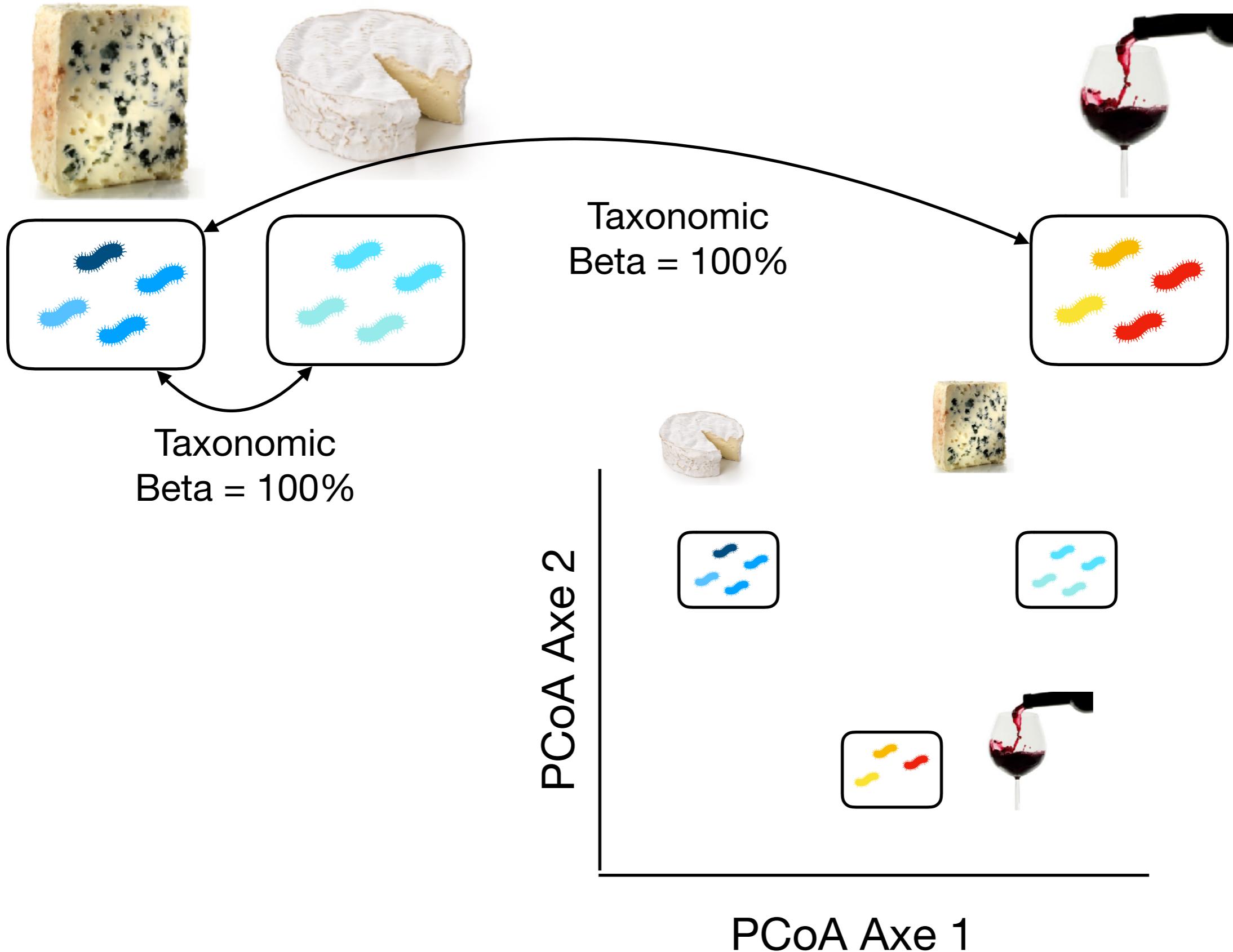
The limits of the OTU approach



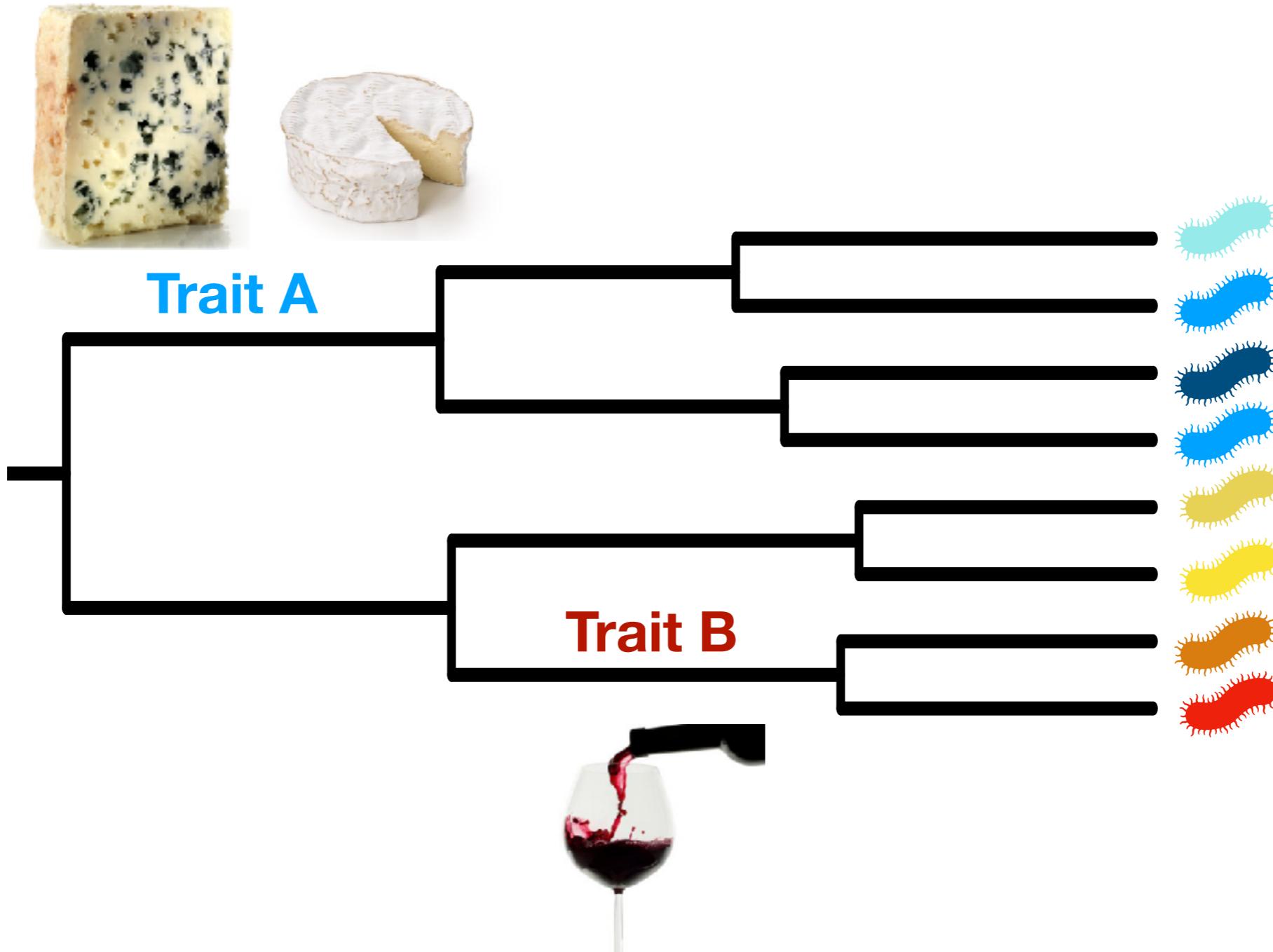
Taxonomic
Beta = 100%



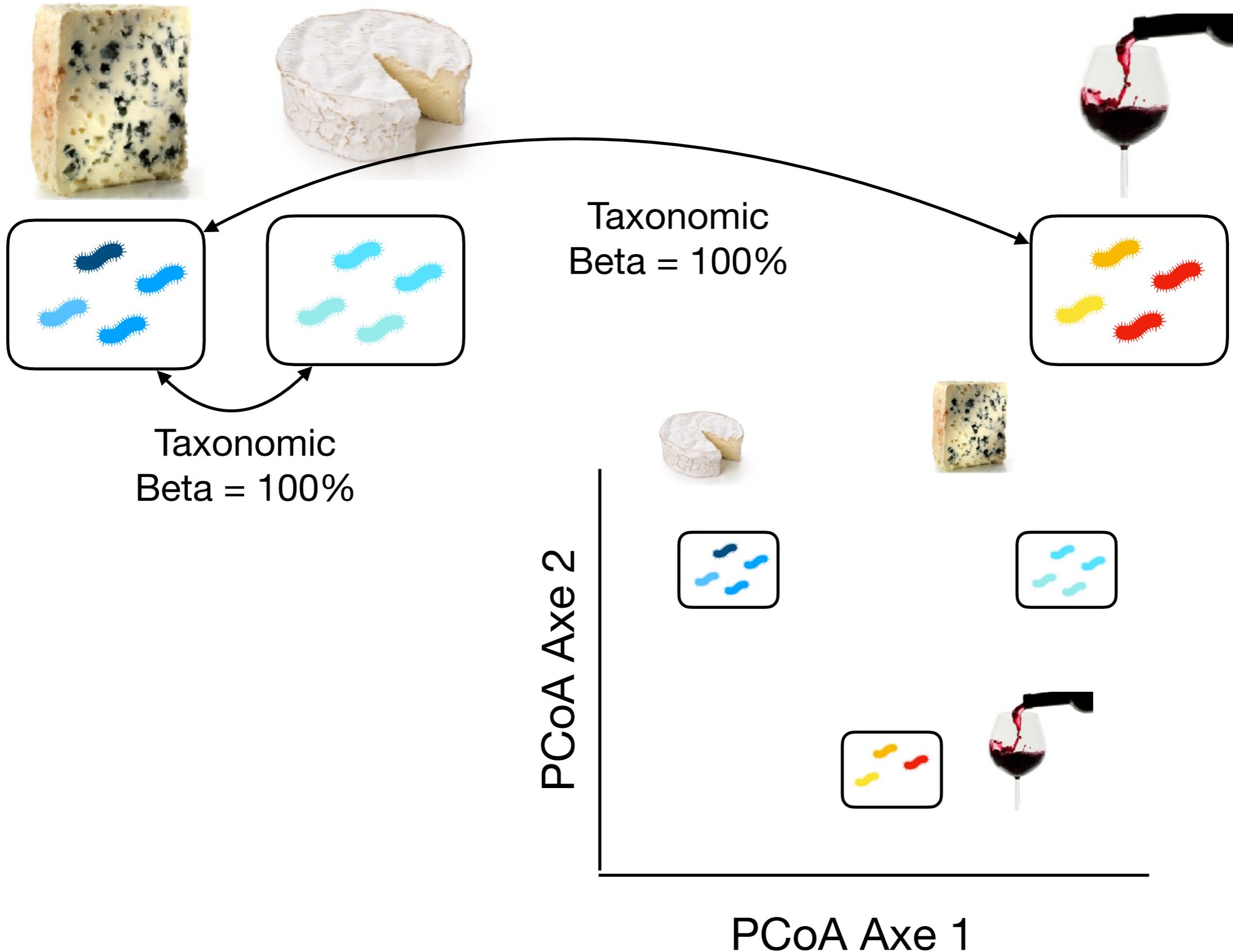
The limits of the OTU approach



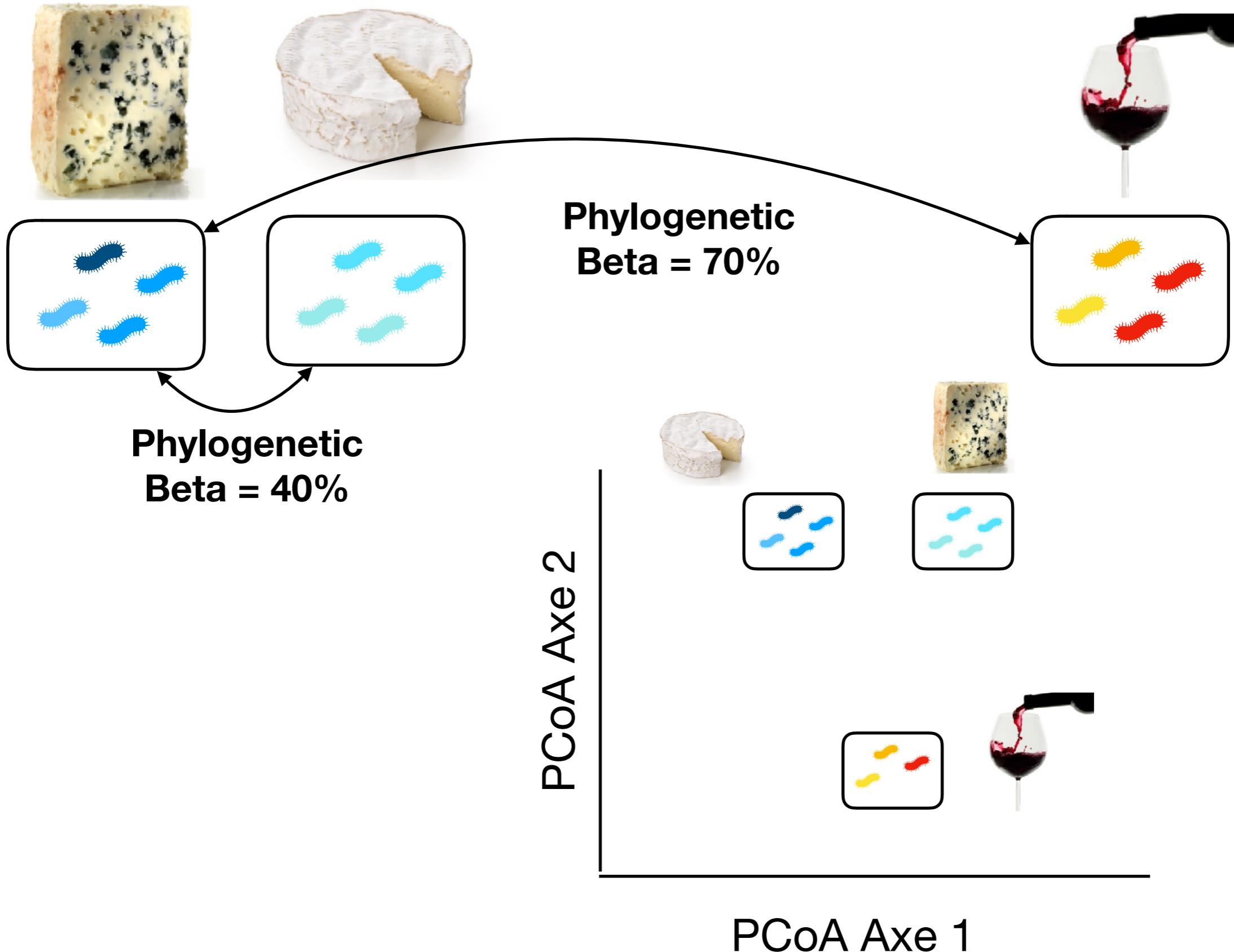
The phylogeny as a proxy for functions



The limits of the OTU approach



The limits of the OTU approach



Objective

**To better describe
the composition of microbiomes
using phylogenetic trees**

Overall structure of the workshop

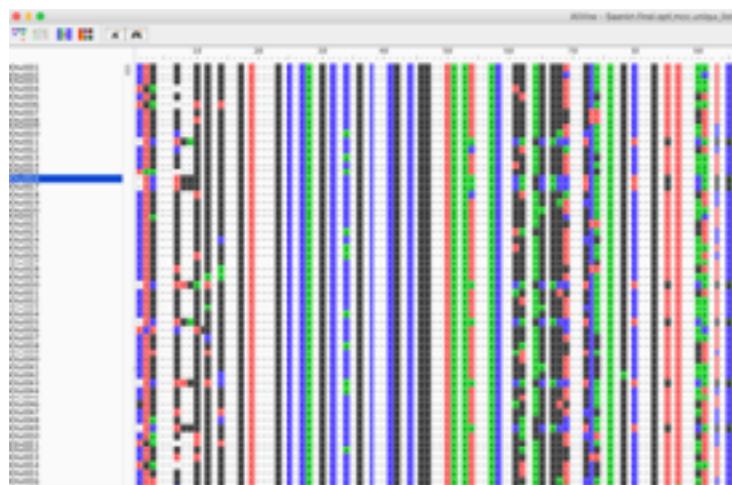
Overall structure of the workshop

0 – The dataset and the softwares

Overall structure of the workshop

0 – The dataset and the softwares

1 – Building a phylogenetic tree



Overall structure of the workshop

0 – The dataset and the softwares

1 – Building a phylogenetic tree

2 – Classical analysis of microbiome

Classical metrics

Bray-Curtis
Unifrac

Classical visualization

PCoA
NMDS

Classical test

PERMANOVA

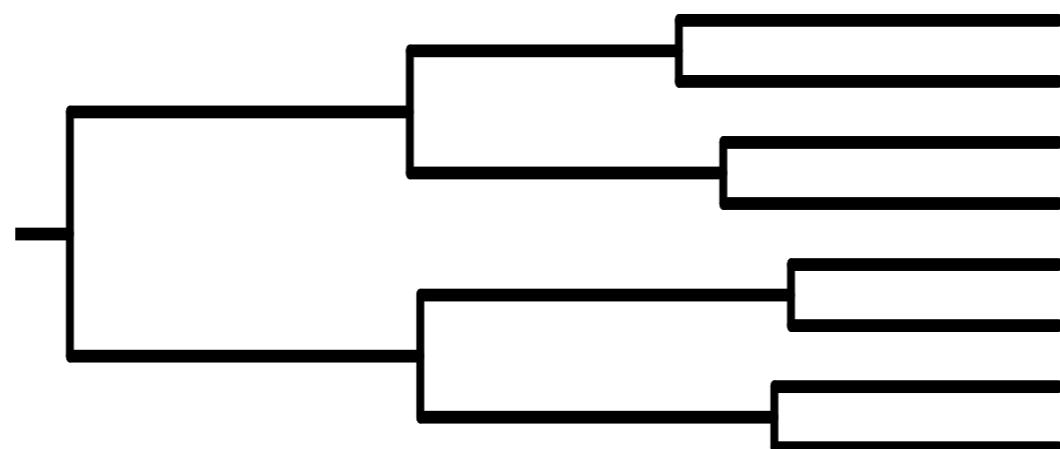
Overall structure of the workshop

0 – The dataset and the softwares

1 – Building a phylogenetic tree

2 – Classical analysis of microbiome

3 – Varying the phylogenetic resolution



Overall structure of the workshop

0 – The dataset and the softwares

1 – Building a phylogenetic tree

2 – Classical analysis of microbiome

3 – Varying the phylogenetic resolution

4 – Exploring the branches of the phylogenetic tree

0 – The dataset and the softwares

Getting started with softwares



Studio



+ Packages

- * tidyverse
 - * ape
 - * seqinR
 - * vegan
- * phyloseq
- * betapart
- * abind
- * Matrix
- * phytools
- * cowplot
- * phylofactor

The screenshot shows the RStudio interface. The left pane displays an R script titled 'Main.R_script.R' with the following content:

```
1 #!/bin/Rscript
2 # Load and install packages
3 # install.packages("tidyverse")
4 # install.packages("phyloseq")
5 # install.packages("ape")
6 # install.packages("seqinR")
7 # install.packages("vegan")
8 # install.packages("cowplot")
9 # install.packages("phylofactor")
10 # Load packages
11 library(tidyverse)
12 library(ape)
13 library(phyloseq)
14 library(vegan)
15 library(phylotools)
16
17 # Set working directory
18 setwd("~/Users/fmazel/Desktop/Research/Phylogenetic/working_directory")
19 setwd("~/Users/fmazel/Desktop/working_directory")
20
21
22 # Load data
23 # Read tree
24 # Plot tree
```

The right pane shows the 'Environment' and 'History' tabs. The 'Environment' tab is empty. The 'History' tab shows the command 'fmazel\$ Desktop/Programmes_Unix/FastTree'.

Below the RStudio interface is a terminal window titled 'fmazel — bash — 80x24' with the following text:

```
Mon Jun 18 09:13:48 on ttys000
Mac:~ fmazel$ Desktop/Programmes_Unix/FastTree
```

FastTree

Where to find the code / the data ?

All material available here:

https://github.com/FloMazel/Microbiome_Phylo_Diversity_Workshop



Organisation of the folders

Practice



39 commits

2 branches

0 releases

2 contributors

Branch: master ▾

New pull request

Create new file

Upload files

Find file

Clone or download ▾

FloMazel	PackagesUpdate	Latest commit 6b33f1e 17 seconds ago
My_outputs	Update_Phlyoseq	6 hours ago
My_outputs_BackUp	Update	7 hours ago
R functions	Update	7 hours ago
data	Merge branch 'pr/3' - update	3 days ago
.DS_Store	Update	3 days ago
.RData	Update	3 days ago
.Rhistory	FastTree_SetUp	5 hours ago
.gitattributes	Initial commit	2 months ago
.gitignore	fixing merge and adding .gitignore	3 days ago
Only_Rscript.R	Add_Links	5 hours ago
README.md	Add_Links	5 hours ago
SetUp.Rmd	PackagesUpdate	17 seconds ago
Workflow.Rmd	PackagesUpdate	17 seconds ago
Workflow.html	Update_Phlyoseq	6 hours ago

Organisation of the folders

Practice



39 commits

2 contributors

1. Download the folder

2. Copy it on your computer

FloMazel PackagesUpdate		
My_outputs	Update_Phlyoseq	6 hours ago
My_outputs_BackUp	Update	7 hours ago
R functions	Update	7 hours ago
data	Merge branch 'pr/3' - update	3 days ago
.DS_Store	Update	3 days ago
.RData	Update	3 days ago
.Rhistory	FastTree_SetUp	5 hours ago
.gitattributes	Initial commit	2 months ago
.gitignore	fixing merge and adding .gitignore	3 days ago
Only_Rscript.R	Add_Links	5 hours ago
README.md	Add_Links	5 hours ago
SetUp.Rmd	PackagesUpdate	17 seconds ago
Workflow.Rmd	PackagesUpdate	17 seconds ago
Workflow.html	Update_Phlyoseq	6 hours ago

Organisation of the folders

Practice



39 commits

2 contributors

1. Download the folder

2. Copy it on your computer

3. The different files

39 commits

2 contributors

Branch: master

New pull

File

Upload files

Find file

Clone or download

FloMazel PackagesUpdate

Latest commit 6b33f1e 17 seconds ago

My_outputs

6 hours ago

My_outputs_BackUp

7 hours ago

R functions

7 hours ago

data

3 days ago

.DS_Store

3 days ago

.RData

3 days ago

.Rhistory

5 hours ago

.gitattributes

2 months ago

.gitignore

3 days ago

Only_Rscript.R

5 hours ago

README.md

5 hours ago

SetUp.Rmd

17 seconds ago

Workflow.Rmd

17 seconds ago

Workflow.html

6 hours ago

(i) Data
(ii) R-scripts and workflow
(iii) Outputs (and backup)

fixing merge and adding .gitignore

Add_Links

Add_Links

PackagesUpdate

PackagesUpdate

Update_Phlyoseq

Organisation of the folders

Practice



Exploring the phylogenetic composition of microbiomes

Some of this code has been adapted from the [ECOSCOPE GitHub account](#)

Table of contents

- [1. Prior to the workshop](#)
- [2. Getting started with R](#)
- [3. Building a microbial phylogenetic tree](#)
- [4. Diversity analysis in R](#)
- [5. PhyloFactor](#)

1. Prior to the workshop

Setup

Please come to the workshop with your laptop set up with the required software and data files as described in our [setup instructions](#).

Background

Please read [Hallam SJ et al. 2017. Sci Data 4: 170156](#) "Monitoring microbial responses to ocean deoxygenation in a model oxygen minimum zone" to learn more about the data used in this workshop. You can also check out this [short video](#) showing how the sampling was done!

Data description

The data in this workshop were collected as part of an on-going oceanographic time series program in Saanich Inlet, a seasonally anoxic fjord on the East coast of Vancouver Island, British Columbia.

2. Getting started with R

Installing and loading packages

At the beginning of every R script, you should have a dedicated space for loading R packages. R packages allow any R user to code reproducible functions and share them with the R community. Packages exist for anything ranging from microbial ecology to complex graphics to multivariate modelling and beyond.

In this workshop, we will use several packages:

Organisation of the folders

Practice

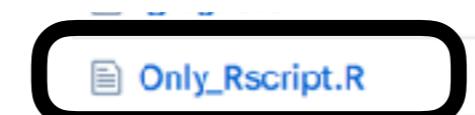


The screenshot shows the RStudio interface with the following components:

- Code Editor:** Displays R code for setting the working directory and importing custom functions. The code includes comments and R Markdown code blocks.
- Console:** Shows the R startup message, license information, and a command to load the seqinr library.
- Help Viewer:** Shows the 'seqinr' package documentation, including sections for 'B' and 'C' containing various functions like as.alignment, as.matrix.alignment, and c2s.

Organisation of the folders

Practice



Only_Rscript.R

```
141 #'
142 ## ---- message=FALSE-----
143 Phylum<-unique(taxonomy$Phylum)
144 Phylum<-subset(Phylum, !(Phylum%in%c("unknown", "Unclassified", "Archaea", "unclassified")) #remove the
145 
146 for (i in Phylum)
147 {
148   taxonomy[[as.character(i)]]$taxonomy$Phylum<-i
149   taxonomy[[as.character(i)]]$Phylum<-i
150 }
151 
152 Constraints<-taxonomy[,c("Bacteria",as.character(Phylum))] #keep only the constraints
head(Constraints)
153 
154 #
155 # Convert to fasta file
156 #
157 ## ---- message=FALSE-----
158 sequences<-list()
159 for (i in 1:dim(Constraints)[1]) sequences[[i]]<-Constraints[i,1]
160 write.fasta(sequences, names=rownames(Constraints), file.out="My_outputs/Phylogenetic_Constraints.fasta", open = "w", nb
161 
162 
163 (Top Level) :
```

Environment

Environment is empty

Files

R. Biological Sequences Retrieval and Analysis

- as.alignment
- as.matrix.alignment
- as.GCNAUCWeb
- as.SeqFastaAA
- as.SeqFastadna
- as.SeqFrag
- autosocket

-- B --

- baselineabit
- tma

-- C --

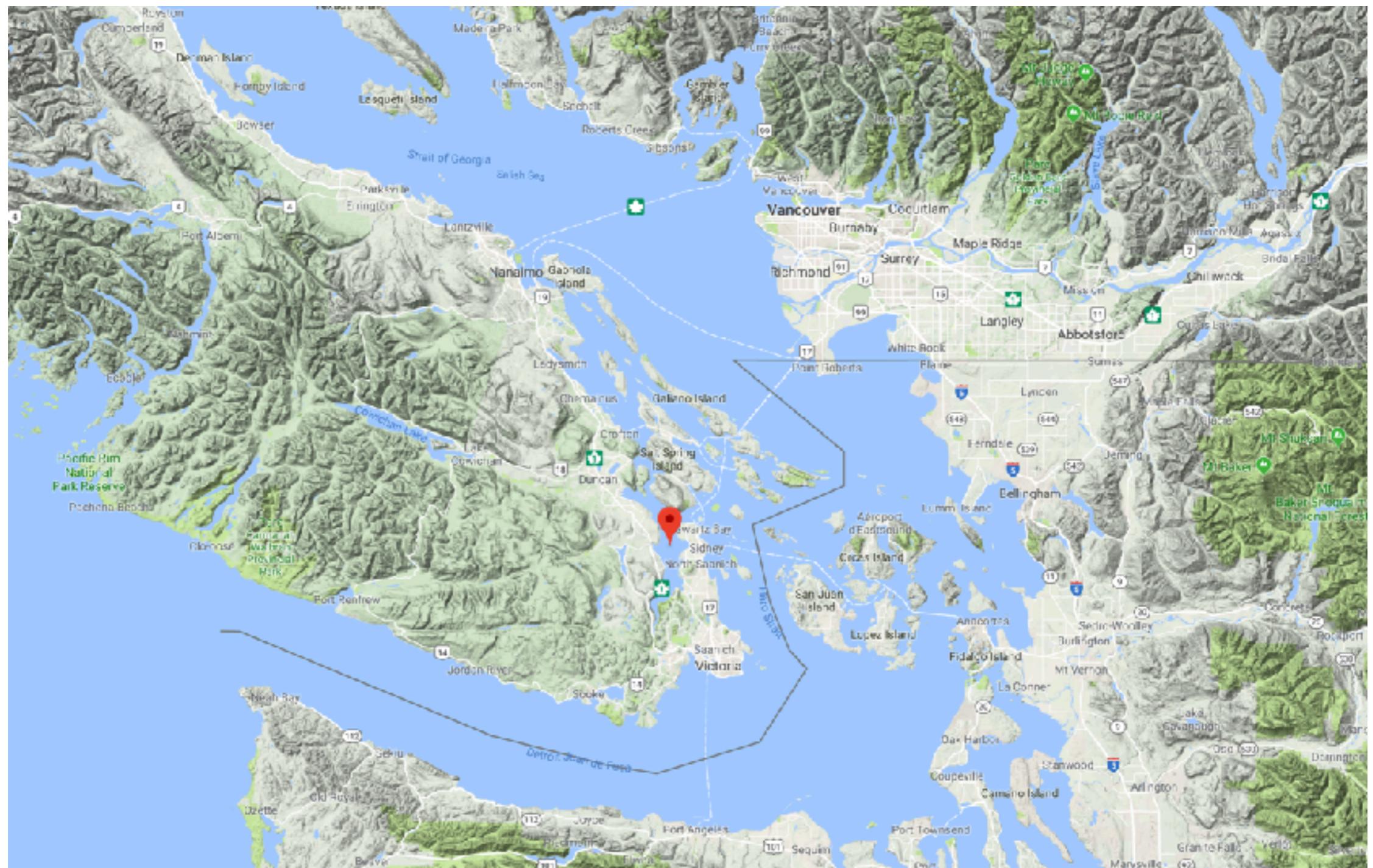
- c2s
- cal
- caltab
- cl
- clnucbase

R is free software and comes with ABSOLUTELY NO WARRANTY.
You are welcome to redistribute it under certain conditions.
Type 'license()' or 'licence()' for distribution details.
Natural language support but running in an English locale
R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'citation()' on how to cite R or R packages in publications.
Type 'demo()' for some demos, 'help()' for on-line help, or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.

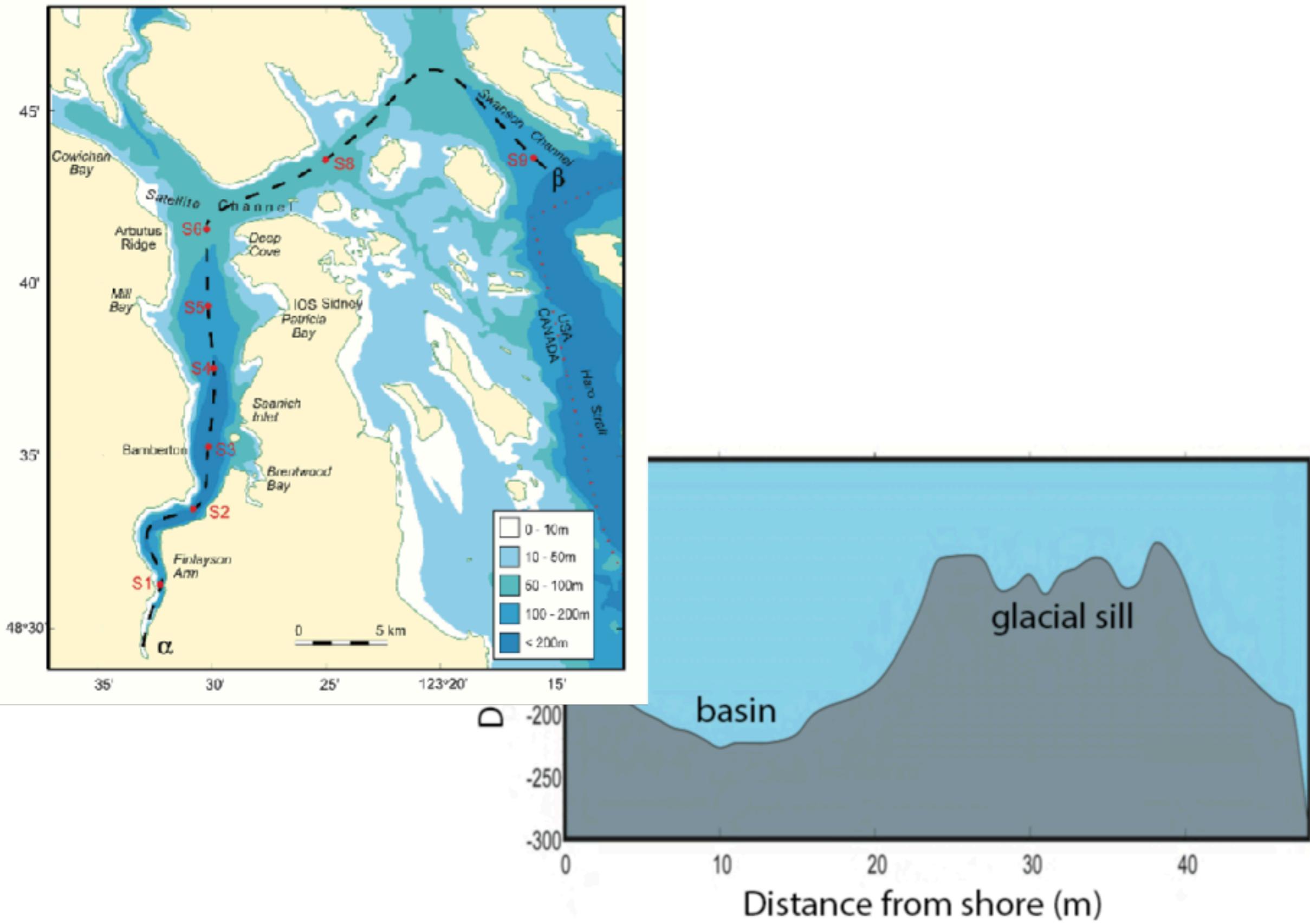
```
> library(seqrinr)
> ?seqinr
>
```

The data

Water columns samples from Saanish inlet

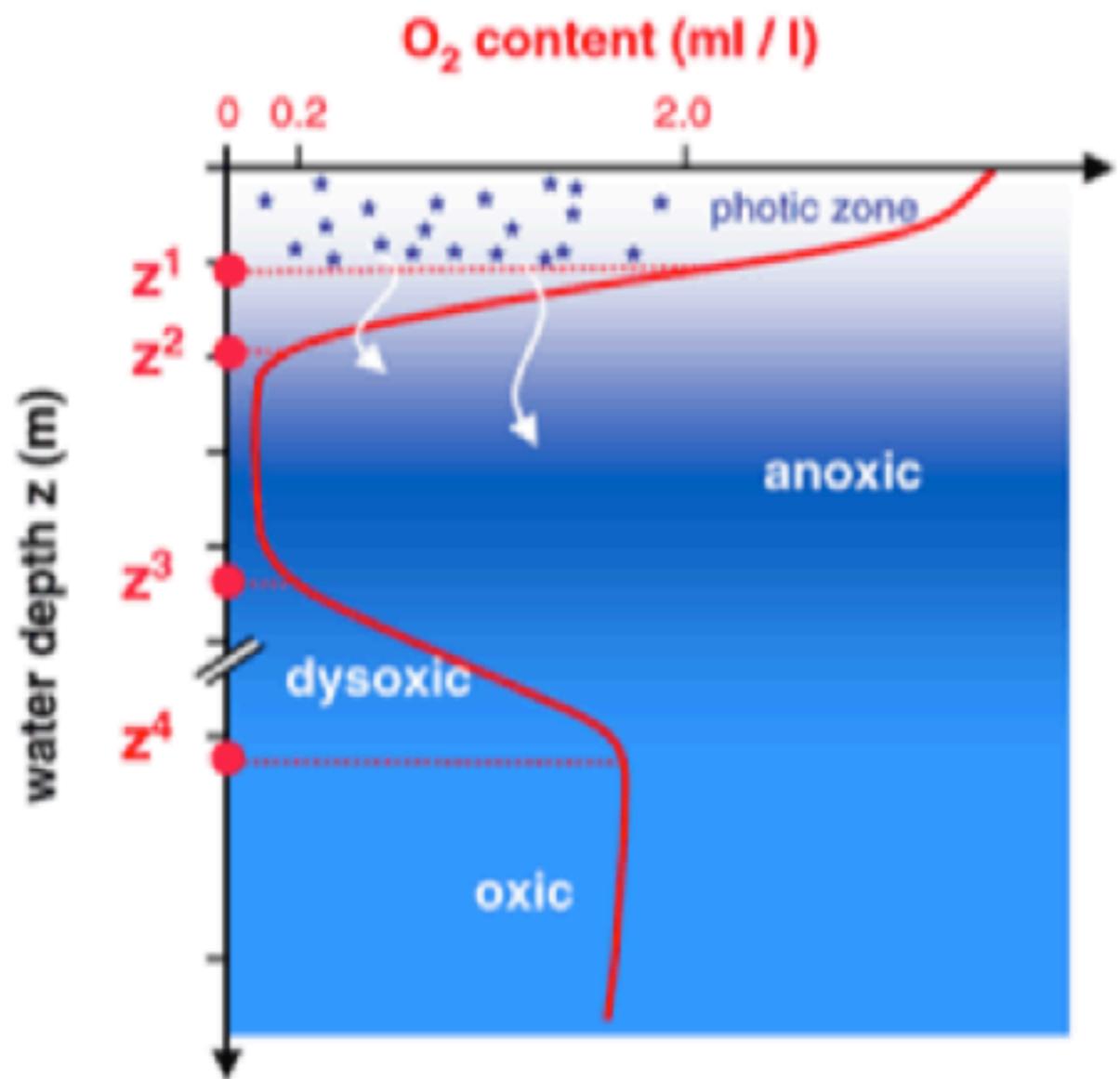


The data



The data

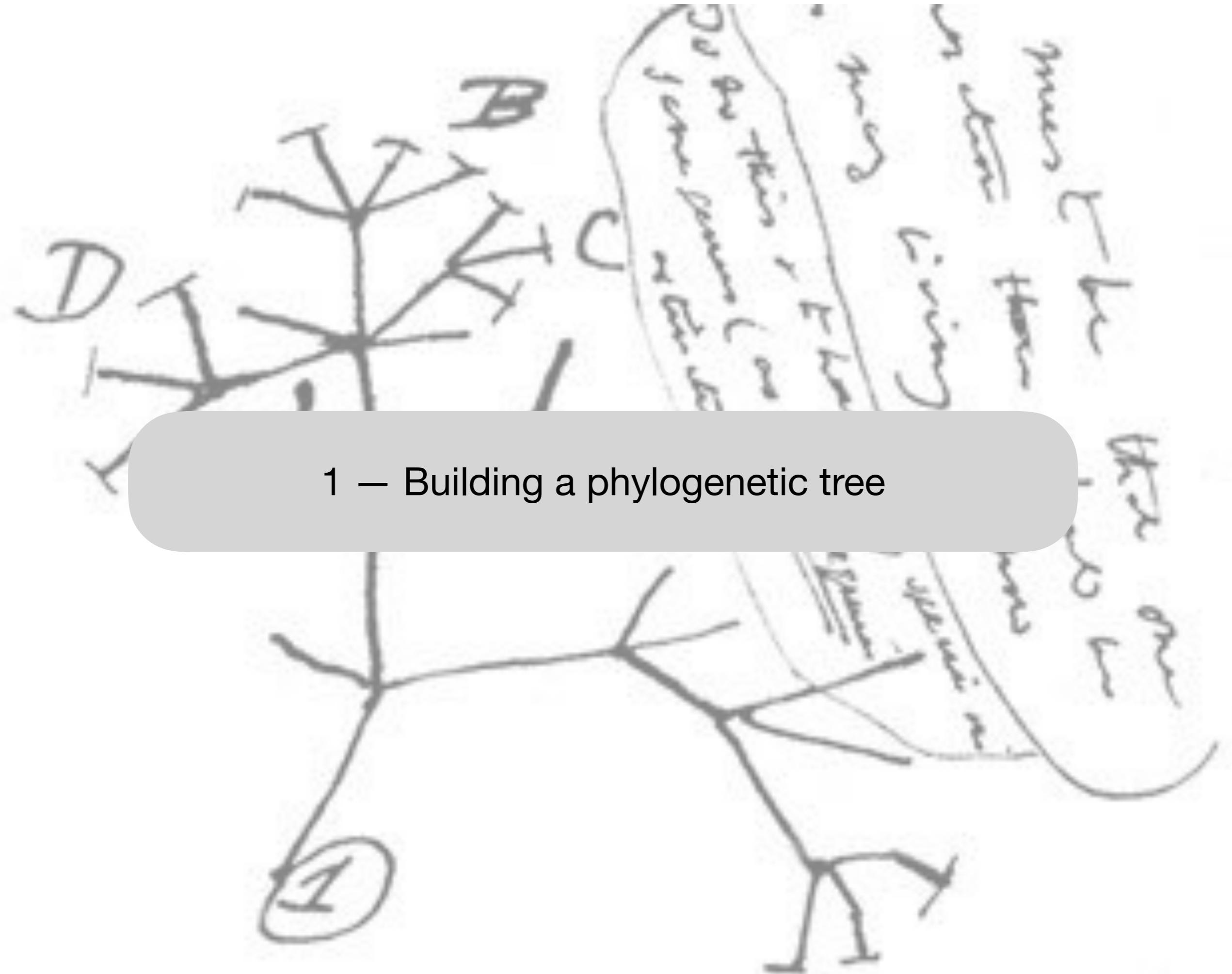
A model for
Oxygen Minimum Zones
(OMZ)



The data

Branch: master ▾	Microbiome_Phylo_Diversity_Workshop / data /	Create new file	Upload files	Find file	History
 FloMazel a					Latest commit da2db6e 16 days ago
..					
 mothur_intermediate_files					16 days ago
 .DS_Store					16 days ago
 Saanich_cruise72_metadata.txt					2 months ago
 Saanich_cruise72_mothur_NJ_tree.RData					2 months ago
 Saanich_cruise72_mothur_OTU_table.shared					2 months ago
 Saanich_cruise72_mothur_OTU_taxonomy.taxonomy					2 months ago
 Saanich_cruise72_mothur_phyloseq.RData					2 months ago
 Saanich_timeseries_geochemical_data.csv		Initial commit			2 months ago
 mothur_pipeline.html		Initial commit			2 months ago

1 – Building a phylogenetic tree

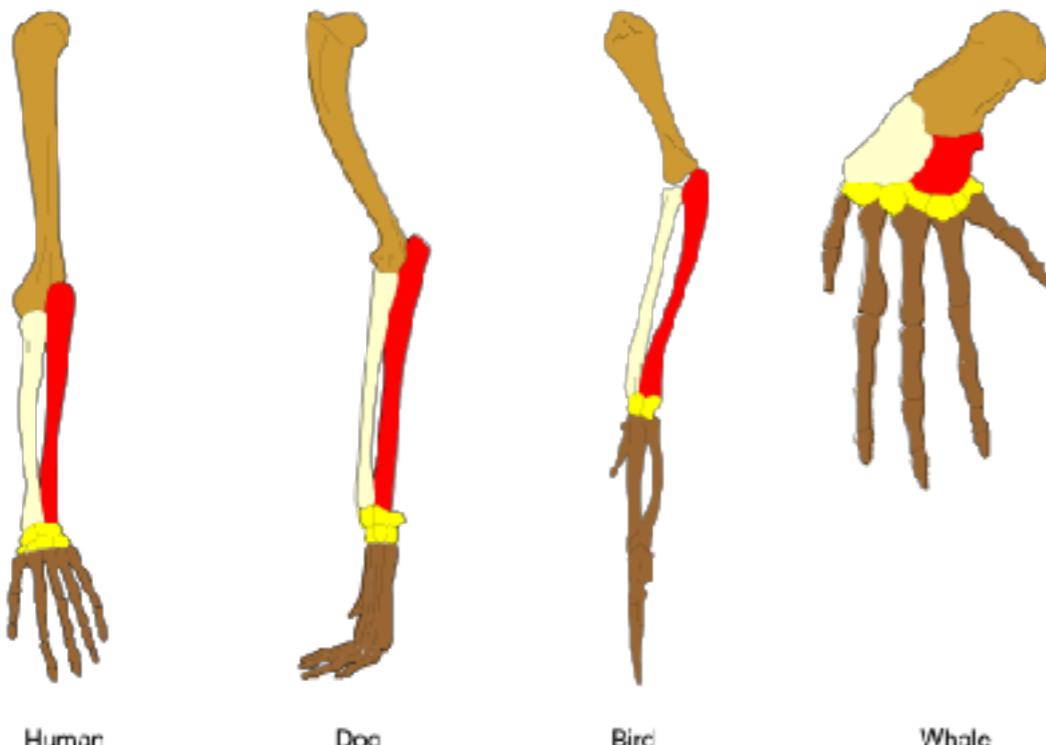


(1) Define homologous characters

(2) Construct a tree

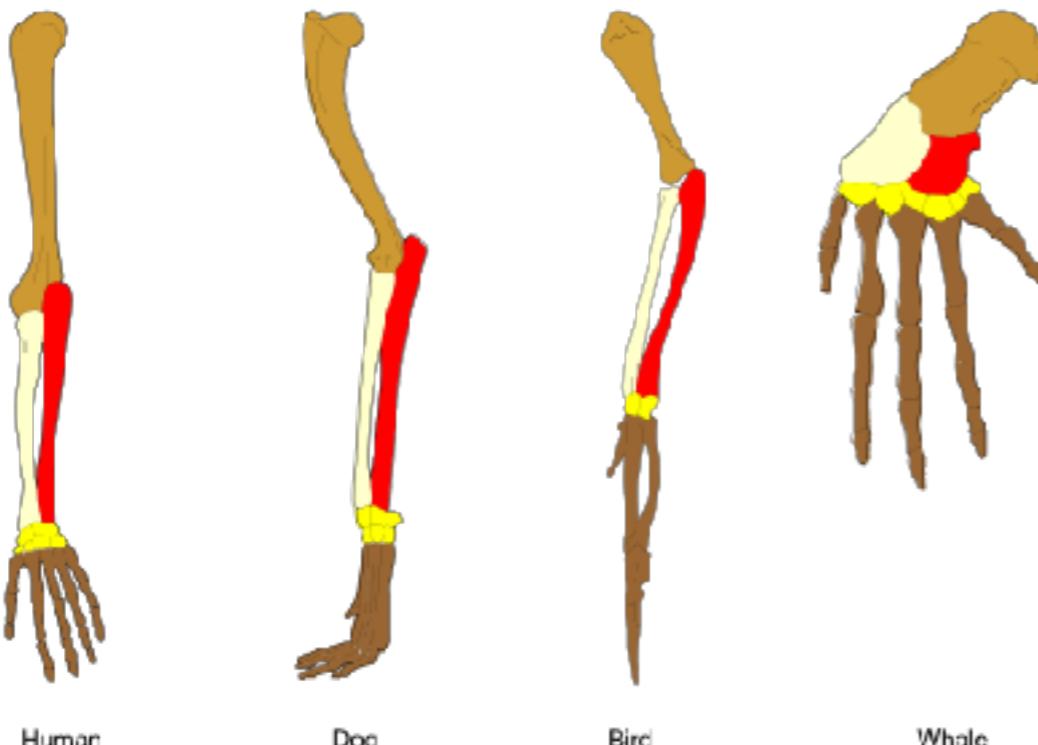
(1) Define homologous characters

Morphological data

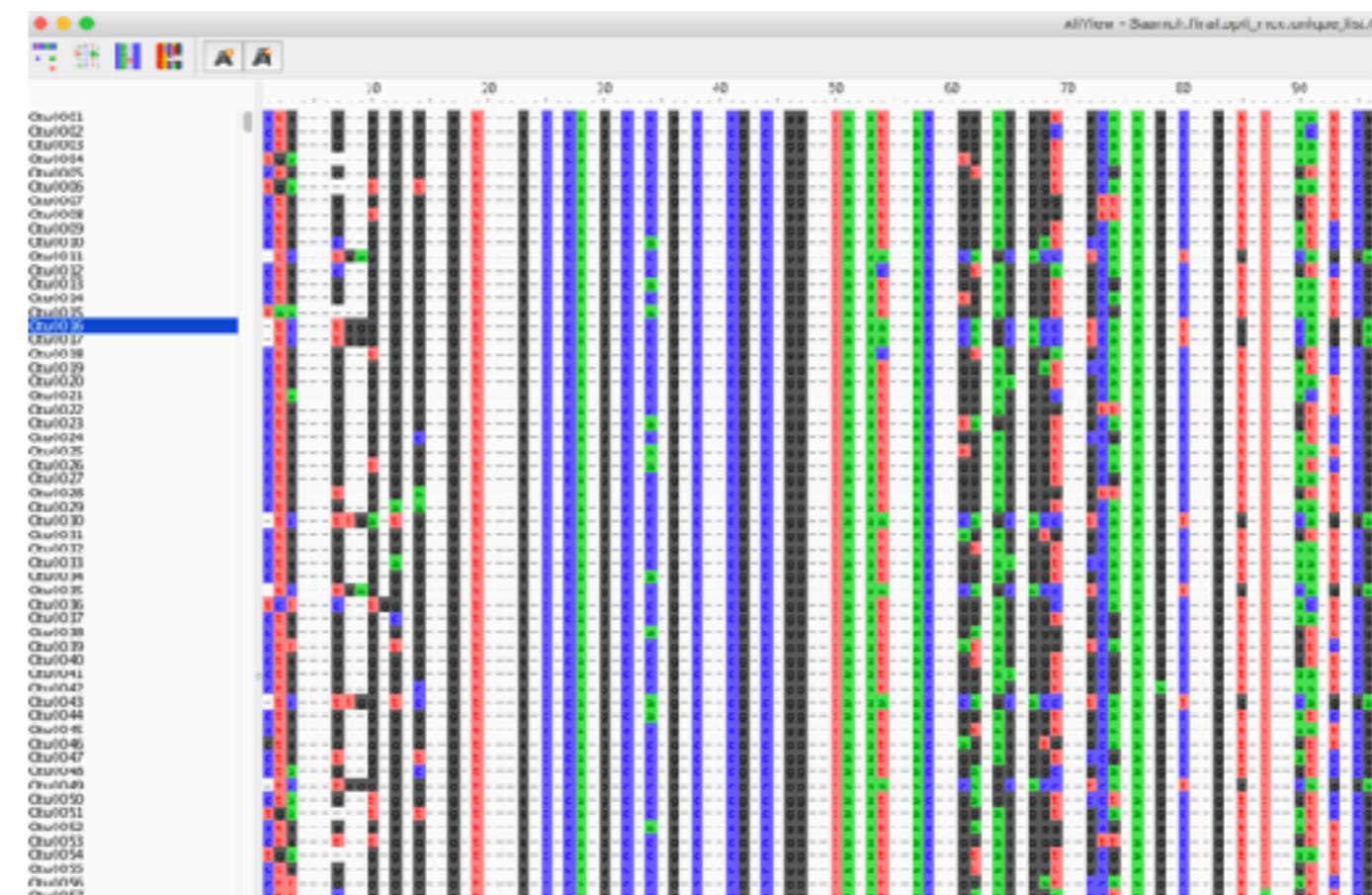


(1) Define homologous characters

Morphological data



Molecular data

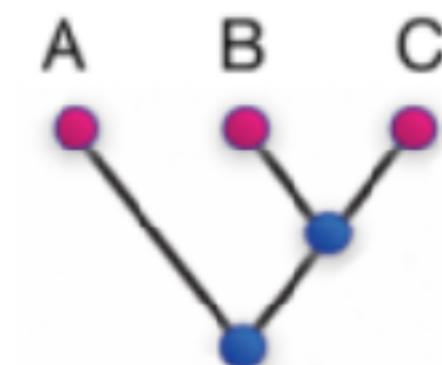
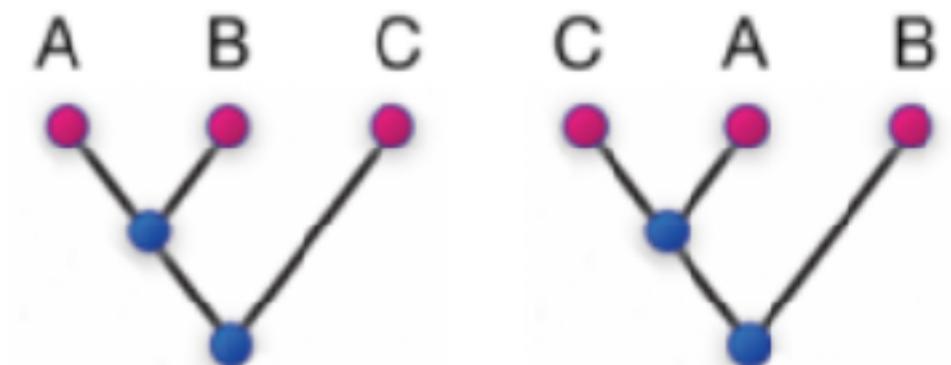
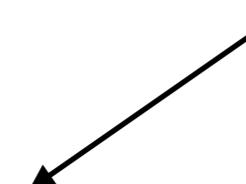


(2) Construct a tree

Idea: Define a criteria to rank trees & explore tree space

Parsimony approach

Maximum likelihood



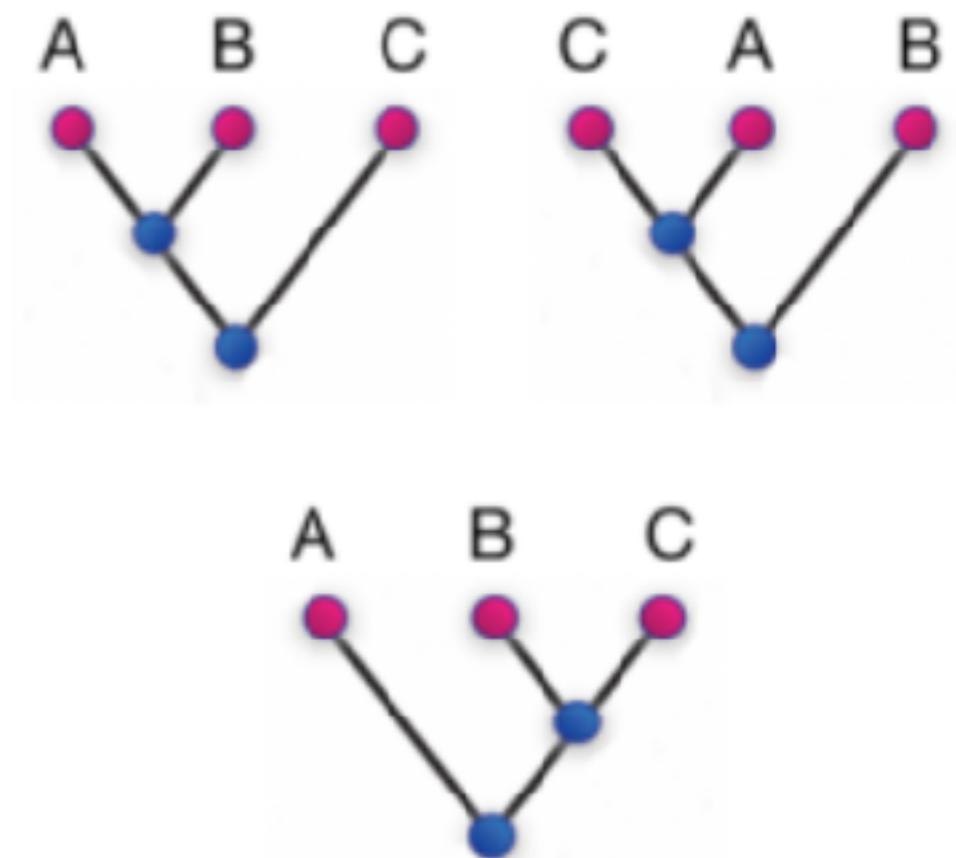
(2) Construct a tree

Idea: Define a criteria to rank trees & explore tree space

Parsimony approach

#ParsimonyGate

Maximum likelihood



(2) Construct a tree

Idea: Define a criteria to rank trees & explore tree space

Parsimony approach

The scenario with the minimal number of changes is the best



Image based on [Taxonomy and phylogeny: Figure 6](#), by Robert Bear et al., CC BY 4.0

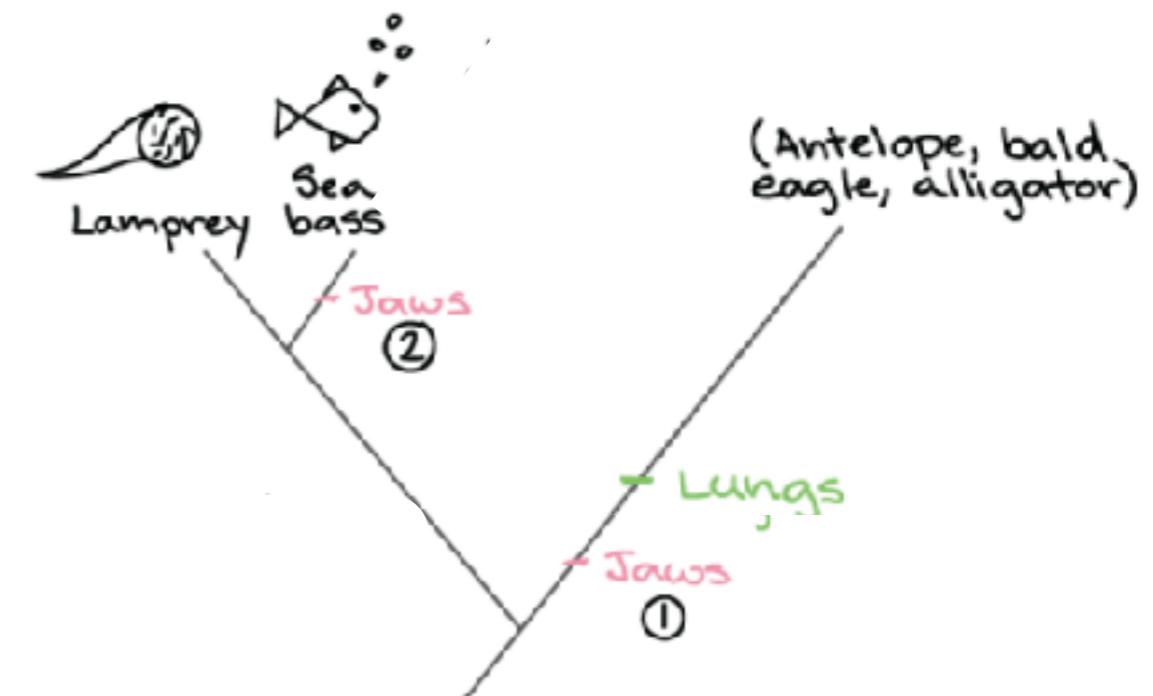
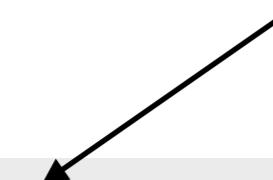


Image based on [Taxonomy and phylogeny: Figure 6](#), by Robert Bear et al., CC BY 4.0

(2) Construct a tree

Idea: Define a criteria to rank trees & explore tree space

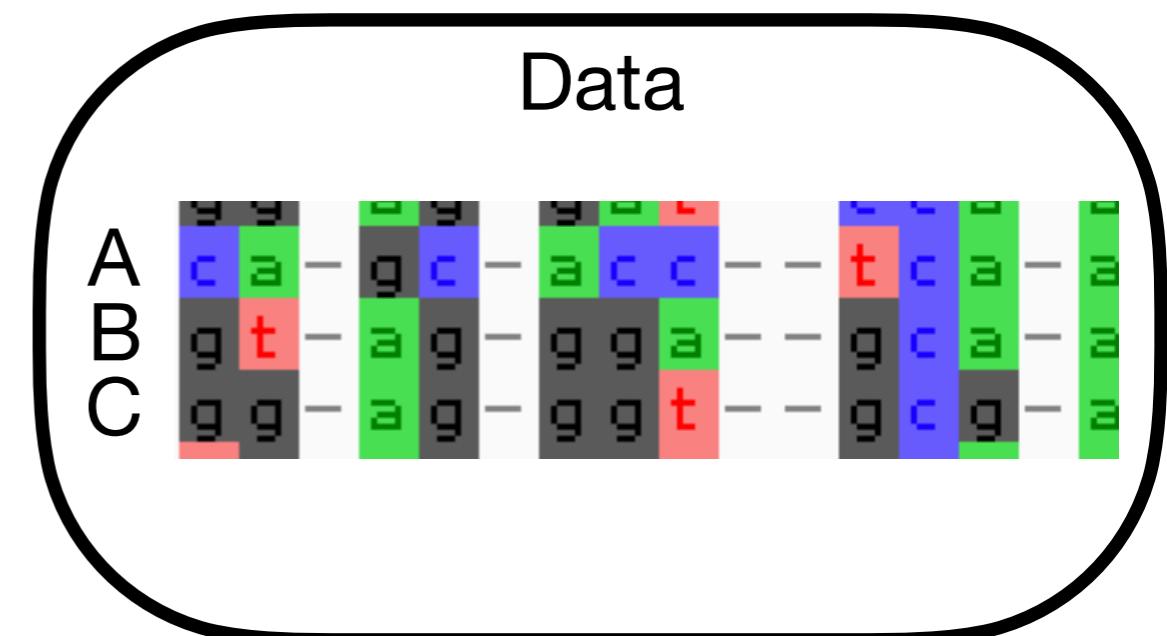
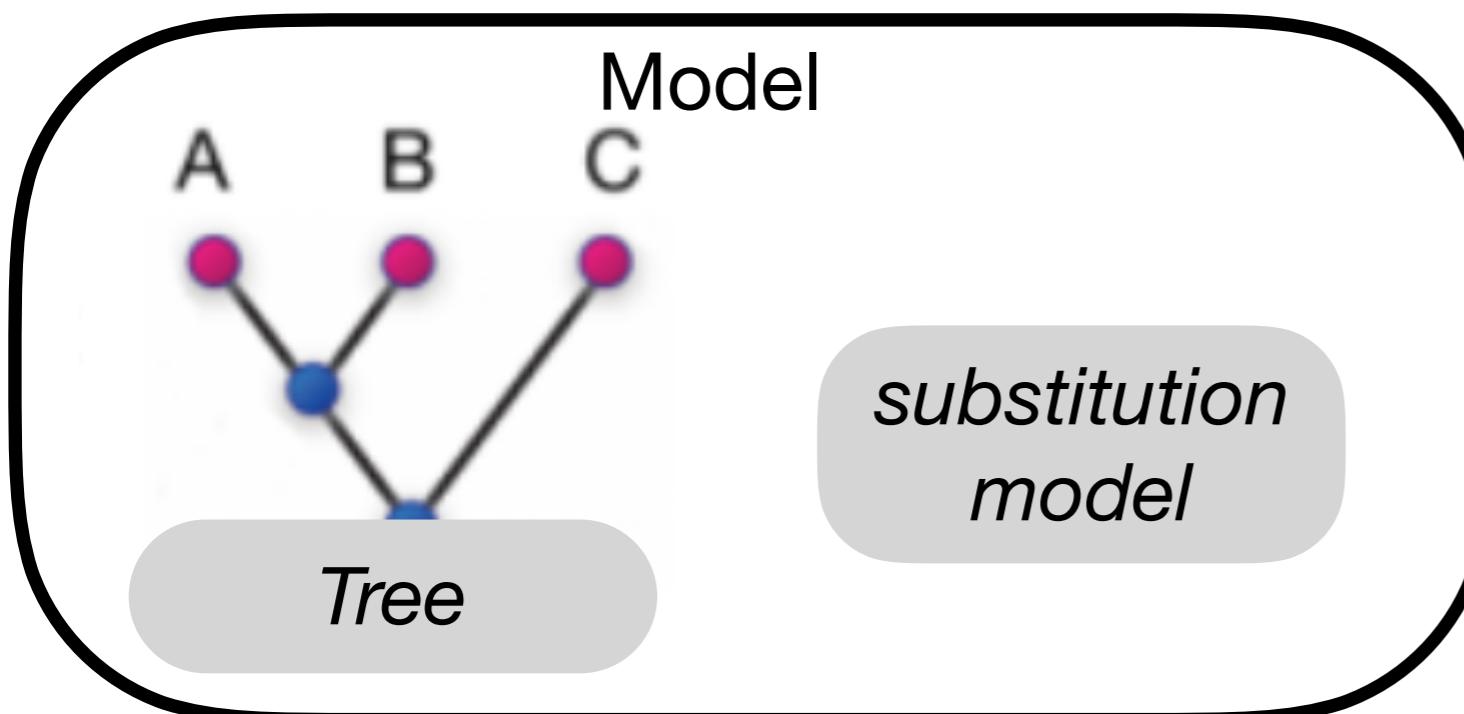
Maximum likelihood



(2) Construct a tree

Idea: Define a criteria to rank trees & explore tree space

Maximum likelihood



Fit: Estimate parameters
+ compute likelihood

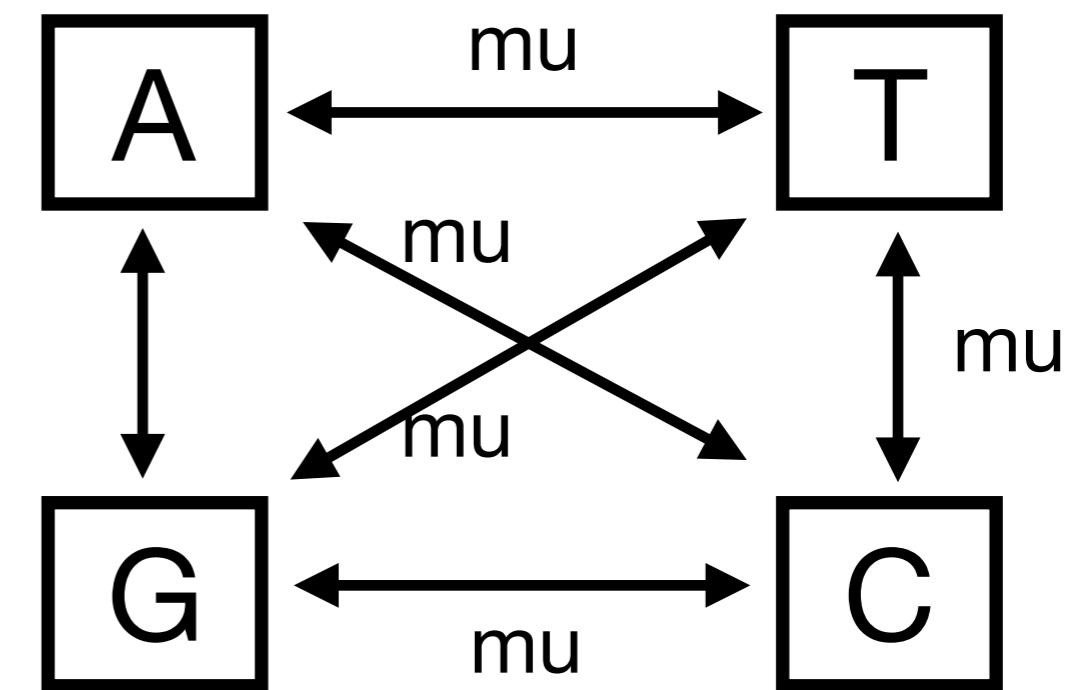
(2) Construct a tree

Idea: Define a criteria to rank trees & explore tree space

Maximum likelihood

Use of a *substitution model*

A lot of different models
exist
(they vary in parameters
richness)



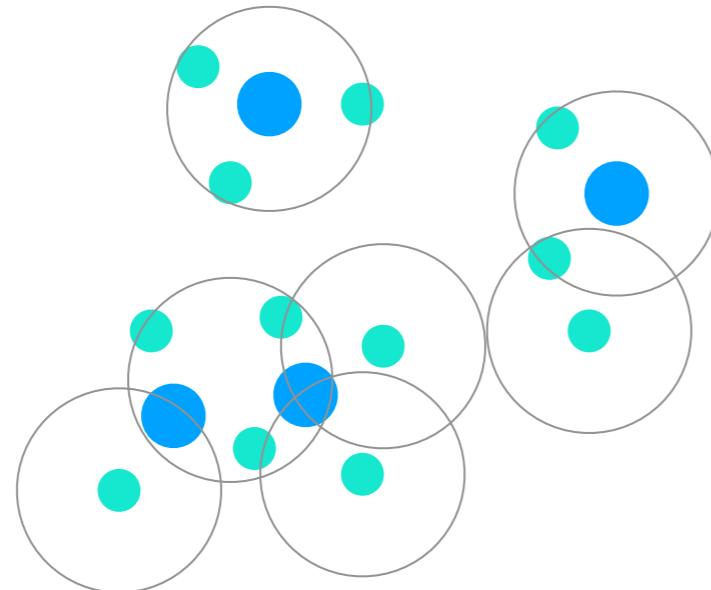
μ = transition rate

New methods for microbiomes:

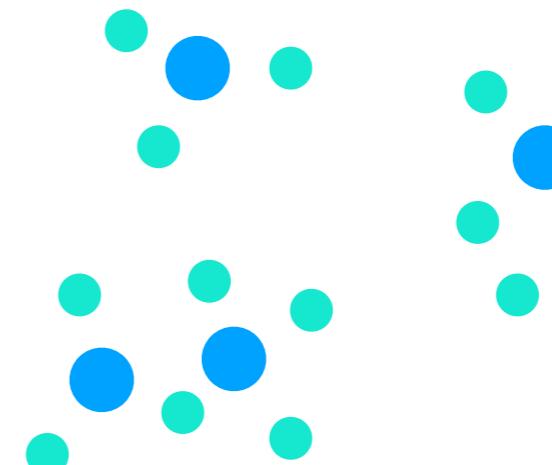
Exact Sequence Variants

- DADA2, Deblur
- Stable IDs
- Less inflation

OTUs



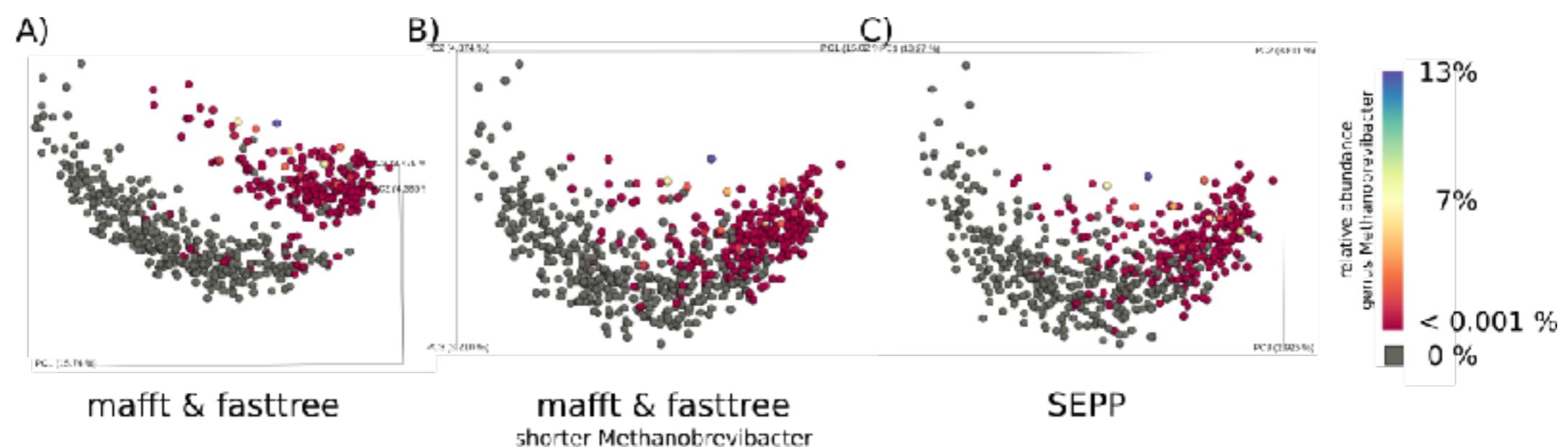
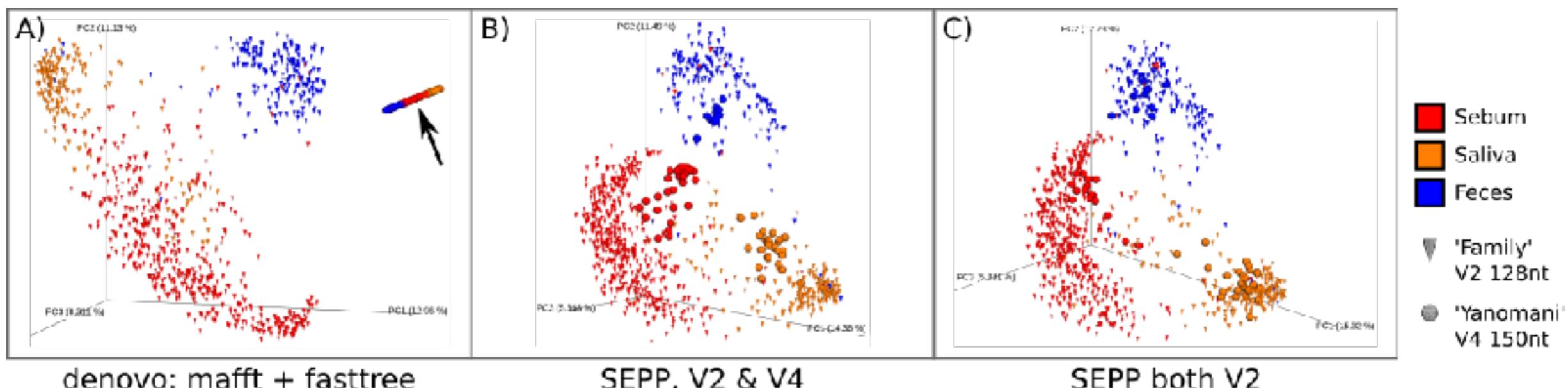
ESVs



New methods for microbiomes:

Phylogenetic Insertion

- q2-fragment-insertion
- fewer artifacts
- more stable meta-analysis



New methods for microbiomes:

Phylogenetic Insertion

Exact Sequence Variants

- both available as Qiime2 plugins
- not covered further here :)

Practice



Build the tree

Part 1 – Constrains

Part 2 – Tree search

Build the tree – Part 1: Constrains

Build the tree – Part 1: Constraints

Rationale

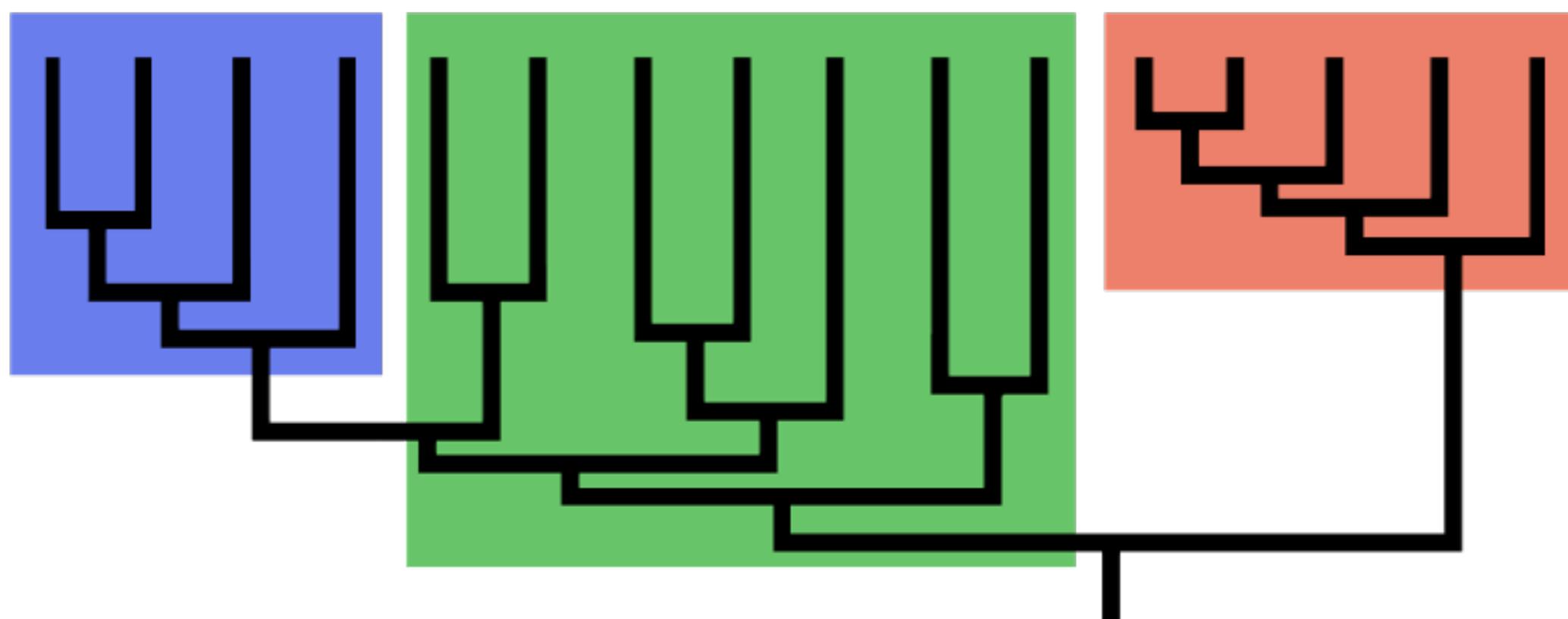
High taxonomic classes (e.g. phylum) have been validated as monophyletic based on longer sequences

Build the tree – Part 1: Constraints

Rationale

High taxonomic classes (e.g. phylum) have been validated as monophyletic based on longer sequences

A **monophyletic** group, or **clade**, is a group of organisms that consists of all the descendants of a common ancestor



Build the tree – Part 1: Constraints

Rationale

High taxonomic classes (e.g. phylum) have been validated as monophyletic based on longer sequences

Idea

Constrain the construction of the tree so that OTUs belonging to the same taxonomic group are grouped in a monophyletic group in the tree

Build the tree – Part 1: Constraints

FastTree -help

Constraints Input
for FastTree

```
fmazels-iMac:~ fmazels$ Desktop/Programmes_Unix/FastTree -help
FastTree 2.1.10 SSE3:
  FastTree protein_alignment > tree
  FastTree < protein_alignment > tree
  FastTree -out tree protein_alignment
  FastTree -nt nucleotide_alignment > tree
  FastTree -nt -gtr < nucleotide_alignment > tree
  FastTree < nucleotide_alignment > tree
FastTree accepts alignments in fasta or phylip interleaved formats

Common options (must be before the alignment file):
  -quiet to suppress reporting information
  -nopr to suppress progress indicator
  -log logfile -- save intermediate trees, settings, and model details
  -fastest -- speed up the neighbor joining phase & reduce memory usage
    (recommended for >50,000 sequences)
  -n <number> to analyze multiple alignments (phylip format only)
    (use for global bootstrap, with seqboot and CompareToBootstrap.pl)
  -nosupport to not compute support values
  -intree newick_file to set the starting tree(s)
  -intree1 newick_file to use this starting tree for all the alignments
    (for faster global bootstrap on huge alignments)
  -pseudo to use pseudocounts (recommended for highly gapped sequences)
  -gtr -- generalized time-reversible model (nucleotide alignments only)
  -lg -- Le-Gascuel 2008 model (amino acid alignments only)
  -wag -- Whelan-And-Goldman 2001 model (amino acid alignments only)
  -quote -- allow spaces and other restricted characters (but not ' ) in
    sequence names and quote names in the output tree (fasta input only;
    FastTree will not be able to read these trees back in)
  -noml to turn off maximum-likelihood
  -neme to turn off minimum-evolution NNIs and SPRs
    (recommended if running additional ML NNIs with -intree)
  -neme -mllen with -intree to optimize branch lengths for a fixed topology
  -cat # to specify the number of rate categories of sites (default 20)
    or -nocat to use constant rates
  -gamma -- after optimizing the tree under the CAT approximation,
    rescale the lengths to optimize the Gamma20 likelihood
  -constraints constraintAlignment to constrain the topology search
    constraintAlignment should have 1s or 0s to indicate splits
  -expert -- see more options
For more information, see http://www.microbesonline.org/fasttree/
```



Build the tree – Part 1: Constrains

Constrains Input
for FastTree

A table OTU*constrained clades
in a fasta file

```
> head(Constrains)
```

	Bacteria	Proteobacteria	Bacteroidetes	Thaumarchaeota	Actinobacteria	Marinimicrobia_(SAR406_clade)	
Otu0001	1	1	0	0	0	0	0
Otu0002	1	1	0	0	0	0	0
Otu0003	1	1	0	0	0	0	0
Otu0004	1	1	0	0	0	0	0
Otu0005	1	1	0	0	0	0	0
Otu0006	1	1	0	0	0	0	0

To construct with the taxonomy file:

```
> head(taxonomy)
```

	Domain	Phylum	Class	Order
Otu0001	Bacteria	Proteobacteria	Gammaproteobacteria	Oceanospirillales
Otu0002	Bacteria	Proteobacteria	Epsilonproteobacteria	Campylobacterales
Otu0003	Bacteria	Proteobacteria	Gammaproteobacteria	Oceanospirillales
Otu0004	Bacteria	Proteobacteria	Gammaproteobacteria	Chromatiales
Otu0005	Bacteria	Proteobacteria	Delta proteobacteria	SAR324_clade(Marine_group_B)
Otu0006	Bacteria	Proteobacteria	Gammaproteobacteria	Oceanospirillales

Build the tree – Part 1: Constraints



Constraints Input for FastTree

Code lines 122 – 185

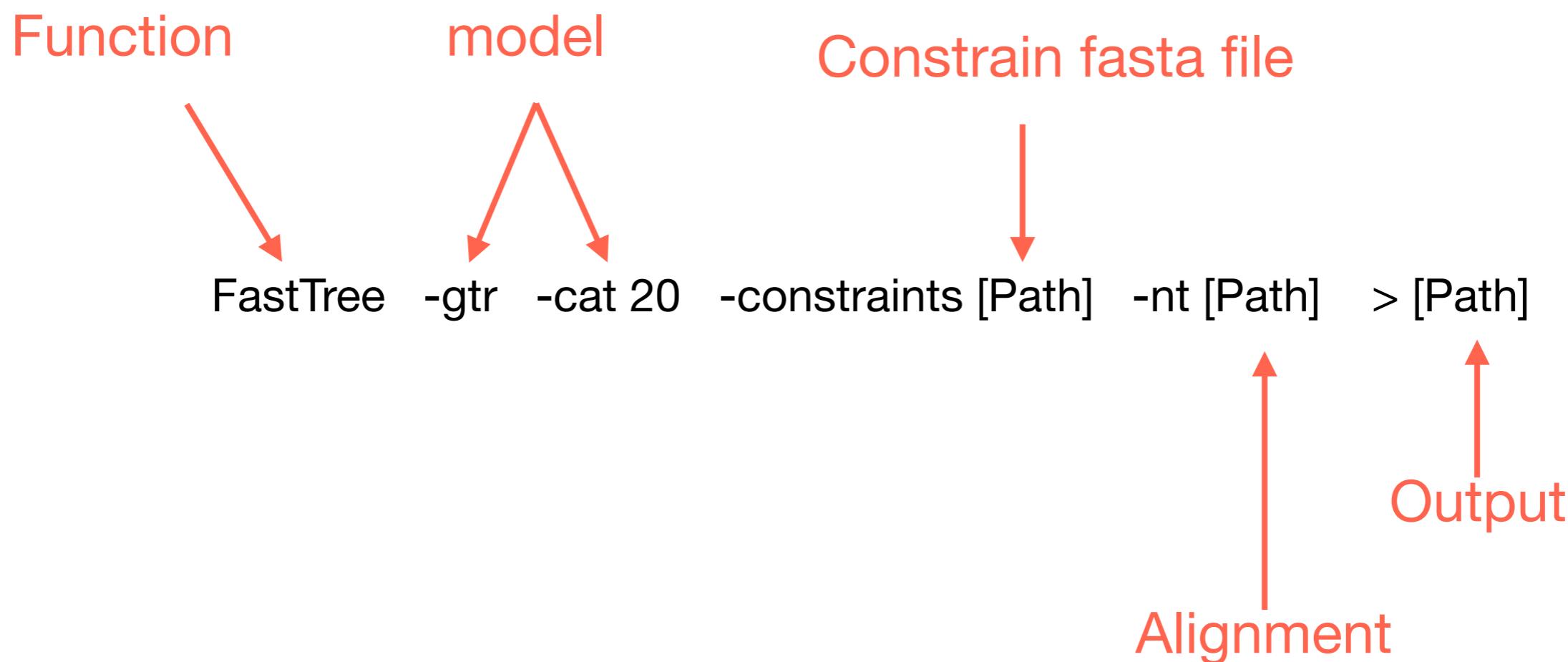


producing the data we are going to use here.

```
115
116 - ## 3.2. Tree Building <a name="Tree-Building"></a>
117
118 Typical read length in microbiome studies are relatively short and have thus contains limited information to reconstruct phylogenetic trees, especially to reconstruct deep branches. To avoid biased phylogenies, we thus constrains deep branches to follow taxonomic classification as it is admitted that large taxonomic clades are monophyletic (informations based on longer sequences). The choice of the constraints is not easy. Here we will constrain tree reconstruction by Domain (Bacteria/Archea) and phylums.
119
120 We will use FastTree to reconstruct the phylogenetic hypotheses. While fastTree is not the best software to reconstruct phylogenies (because it makes a lot of approximations), it has the huge avantage to be fast, which is often critical in microbial species, as there are a very high number of "species".
121
122 - #### Topological constraints
123
124 Load the taxonomic file
125
126 - ```{r, message=FALSE}
127 taxonomy.raw = read.table("data/Saanich_cruise72_mothur_OTU_taxonomy.taxonony", sep="\t", header=TRUE, row.names=1)
128 taxonomy= taxonomy.raw %>%
129   select(-Size) %>%
130   separate(Taxonomy, c("Domain", "Phylum", "Class", "Order", "Family", "Genus", "Species"), sep=";")
```

Build the tree – Part 2 : Run FasTree

On the Terminal



Build the tree – Part 2 : Run FasTree

On the Terminal

Code lines 185 – 211

Build the tree – Part 3 : Plot the tree

Code lines 211 – 260

Load the trees in R

```
> tree=read.tree()
```

Explore the tree file

Colors of the tips

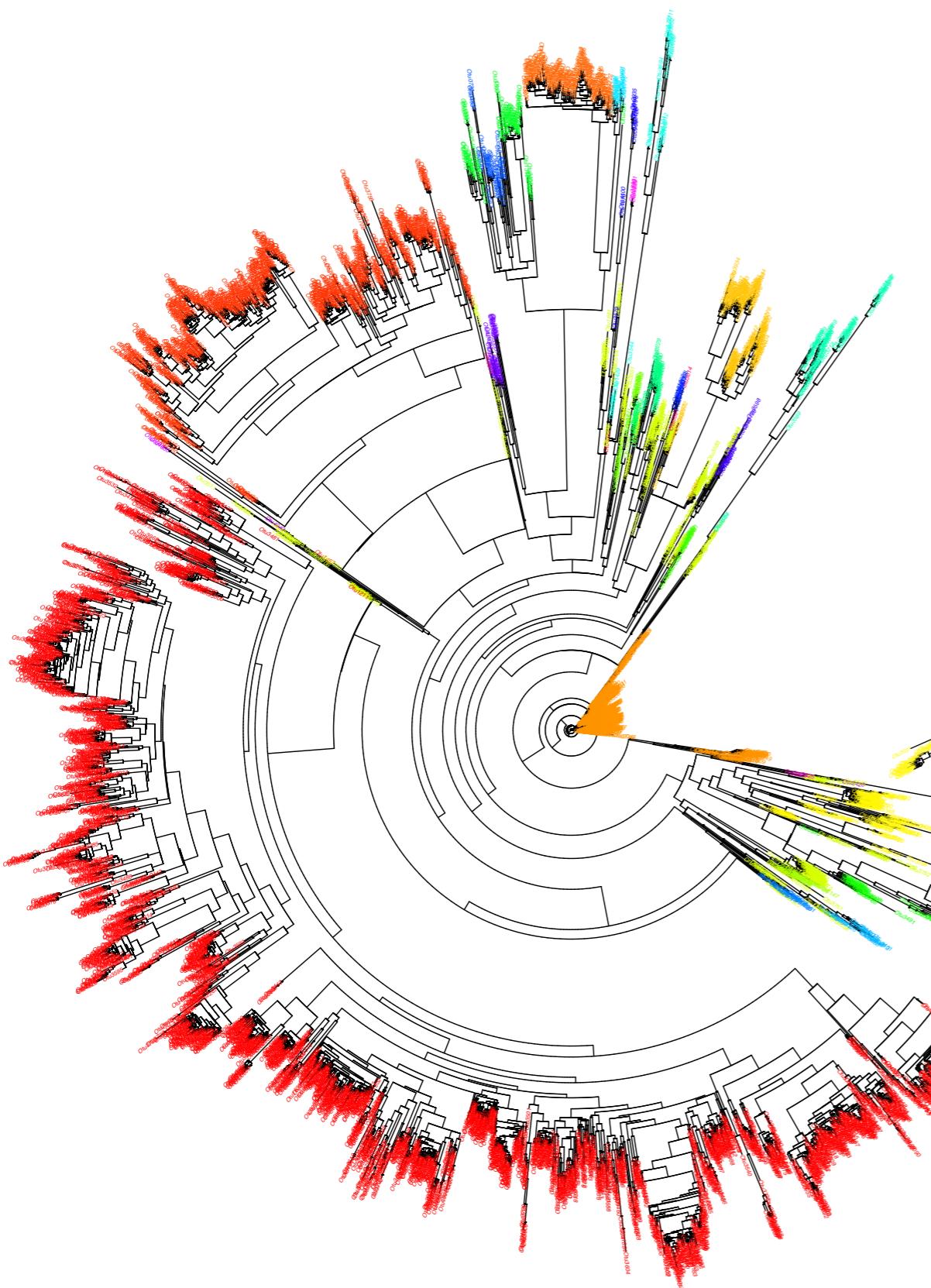
```
> colors=
```

Plot the trees in R

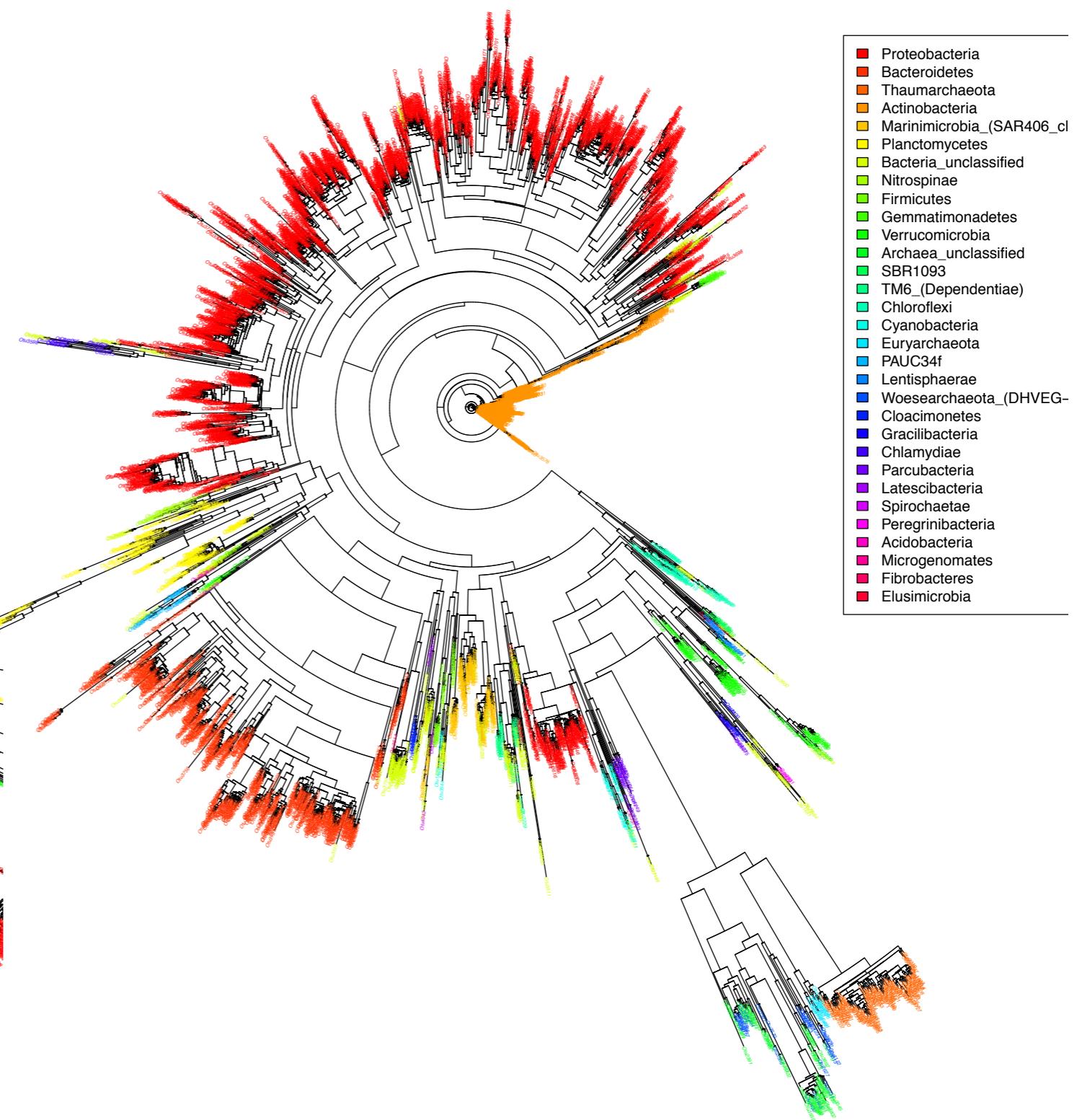
```
> pdf()  
> plot(tree, col=colors)  
> dev.off()
```

Build the tree – Part 3 : Plot the tree

With constraints



Without constraints



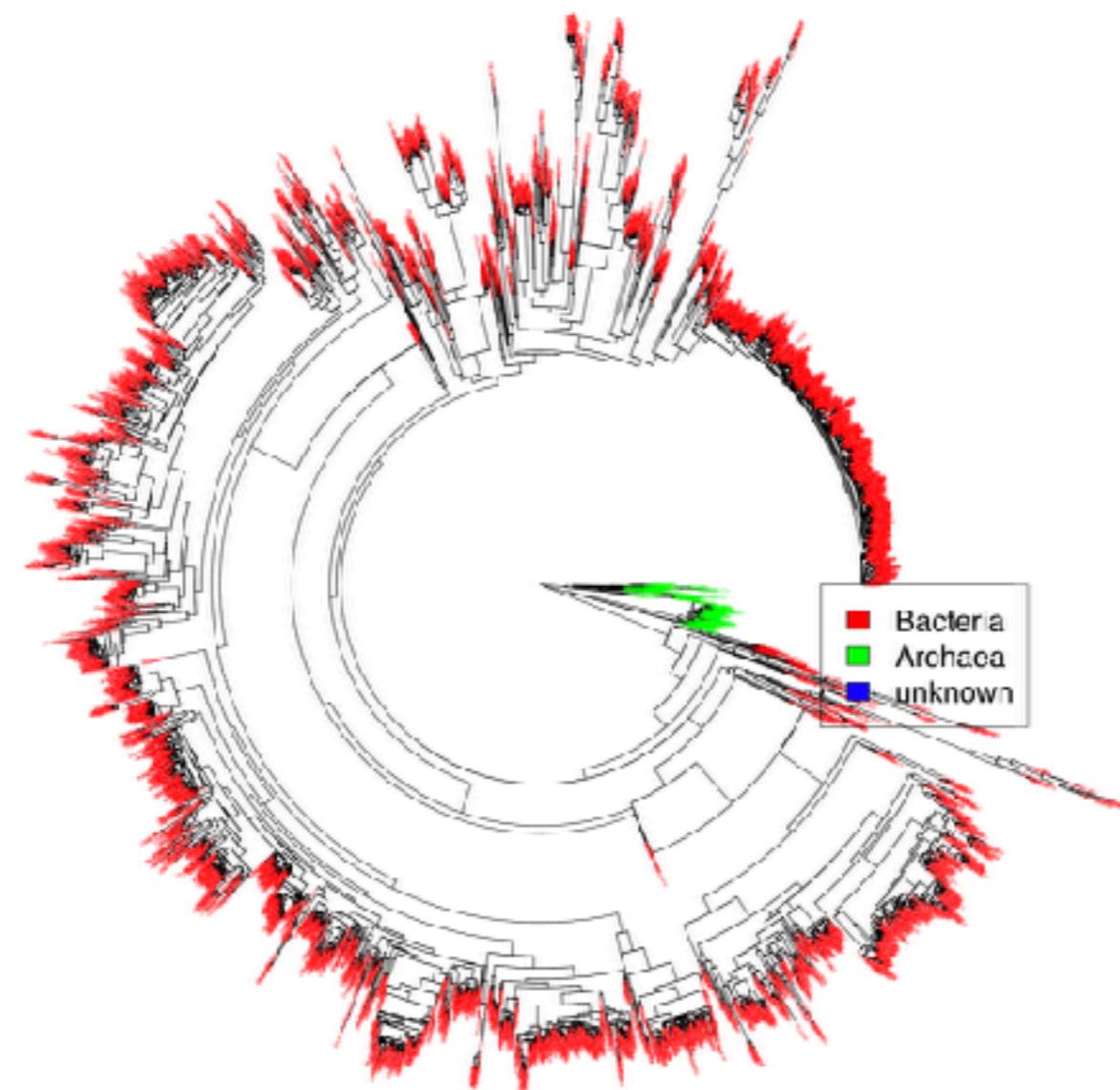
- Proteobacteria
- Bacteroidetes
- Thaumarchaeota
- Actinobacteria
- Marinimicrobia (SAR406 cl)
- Planctomycetes
- Bacteria unclassified
- Nitrospinae
- Firmicutes
- Gemmatimonadetes
- Verrucomicrobia
- Archaea unclassified
- SBR1093
- TM6 (Dependentiae)
- Chloroflexi
- Cyanobacteria
- Euryarchaeota
- PAUC34f
- Lentisphaerae
- Woesearchaeota (DHVEG-)
- Cloacimonetes
- Gracilibacteria
- Chlamydiae
- Parcubacteria
- Latescibacterae
- Spirochaetae
- Peregrinibacteria
- Acidobacteria
- Acidogenomates
- Fibrobacteres
- Elusimicrobia

Build the tree – Part 3 : Root the tree

Code lines 260 – 295

Root the tree in R

> `tree=root()`



Overall structure of the workshop

1 – Building a phylogenetic tree

2 – Classical analysis of microbiome

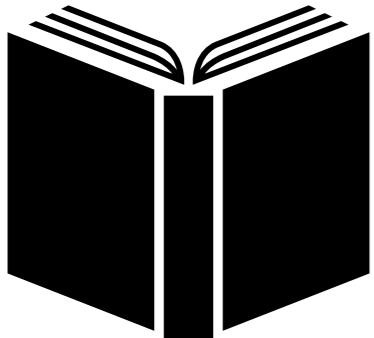
3 – Varying the phylogenetic resolution

4 – Exploring the branches of the phylogenetic tree

2 – Classical analysis of microbiome compositions

“Classical analysis” of the microbiota composition

Theory



Beta-Diversity metric

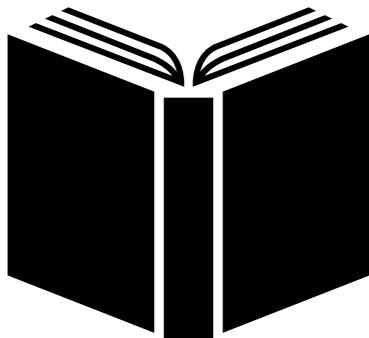
Visualization

Statistical test



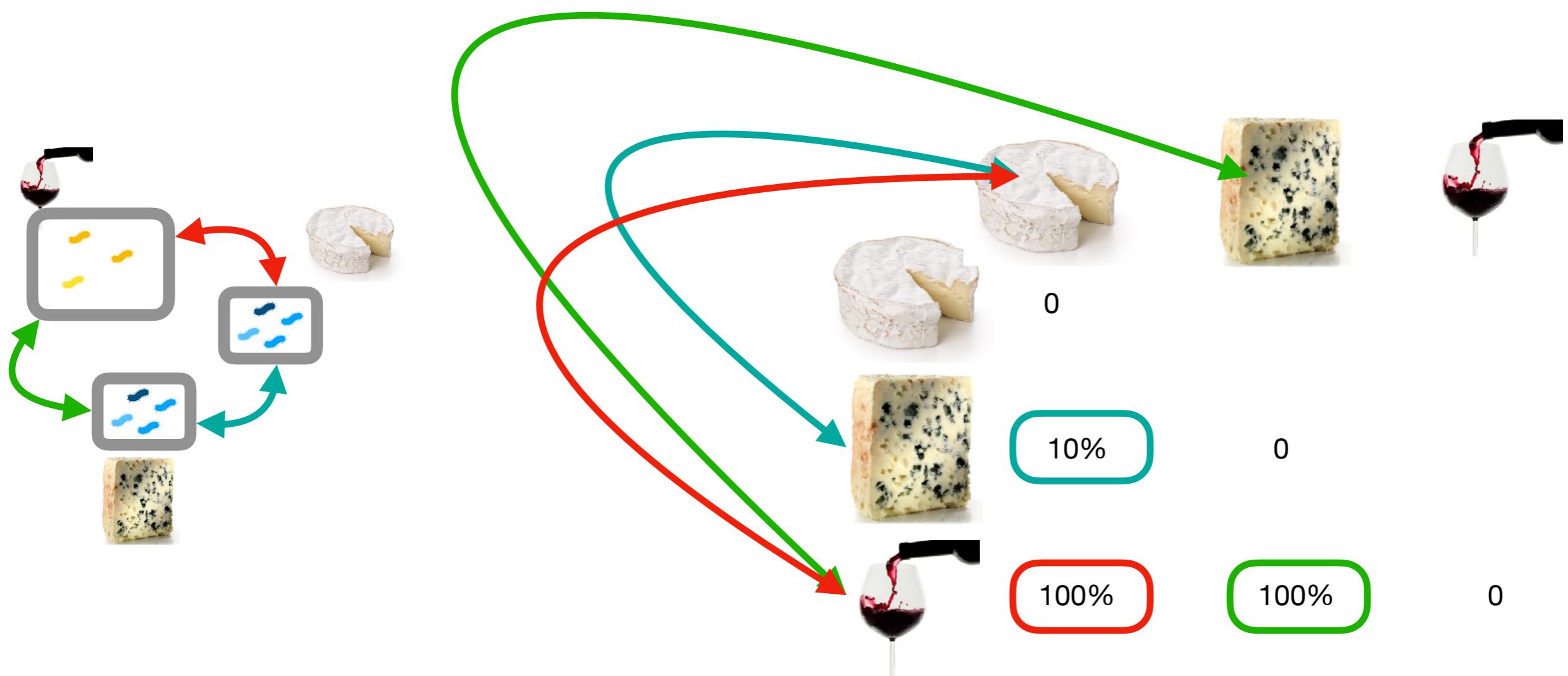
“Classical analysis” of the microbiota composition

Theory



Beta-Diversity metric

Distance matrix



“Classical analysis” of the microbiota composition

Theory



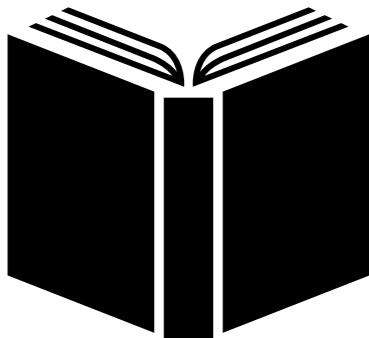
Beta-Diversity metric

Non – Phylogenetic

Presence/Absence

“Classical analysis” of the microbiota composition

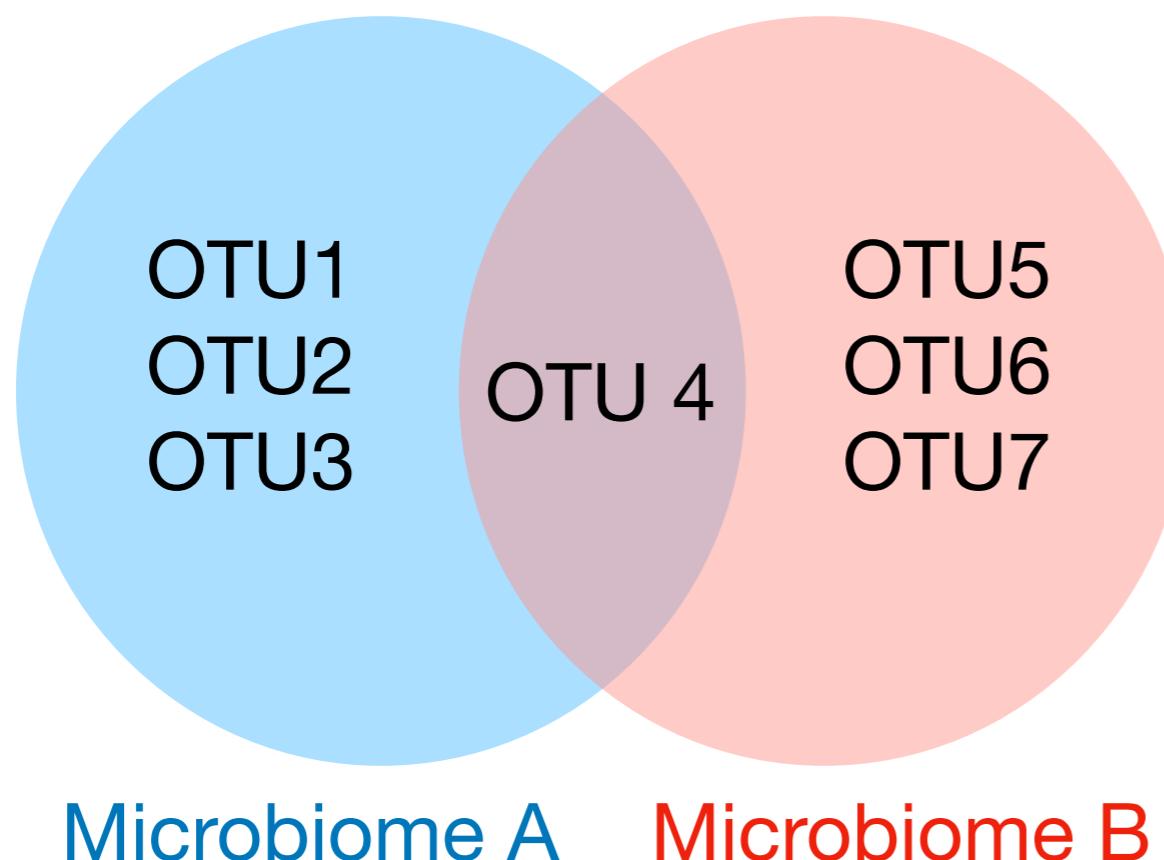
Theory



Beta-Diversity metric

How to quantify similarity?
Example of Jaccard index

$$J(A, B) = 1 - \frac{A \cap B}{A \cup B}$$

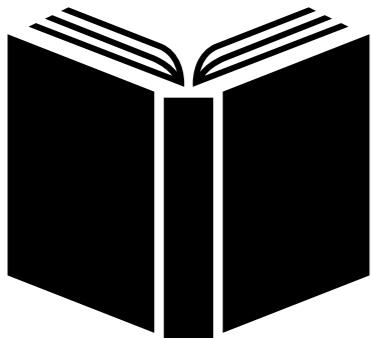


$$A \cap B = 1 \text{ (OTU4)}$$

$$A \cup B = 7 \text{ (OTU1-7)}$$

“Classical analysis” of the microbiota composition

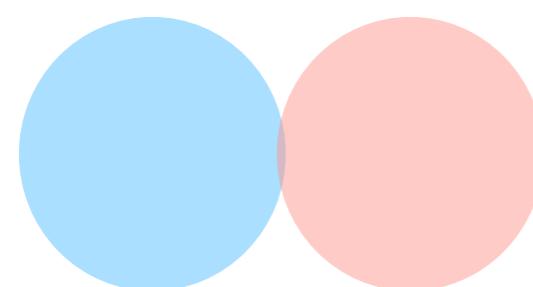
Theory



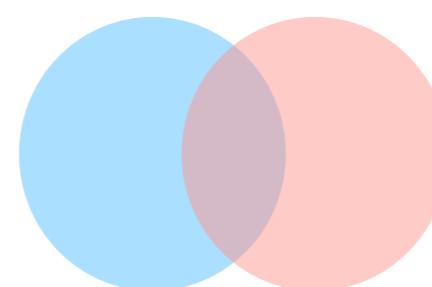
Beta-Diversity metric

How to quantify similarity?
Example of Jaccard index

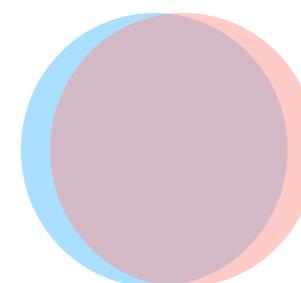
$$J(A, B) = 1 - \frac{A \cap B}{A \cup B}$$



$J=1$



$0 < J < 1$



$J=0$

“Classical analysis” of the microbiota composition

Theory



Beta-Diversity metric

Non – Phylogenetic

Presence/Absence

Jaccard
Sorensen

“Classical analysis” of the microbiota composition

Theory



Beta-Diversity metric

Non – Phylogenetic

**Jaccard
Sorensen**

Presence/Absence

Abundances

“Classical analysis” of the microbiota composition

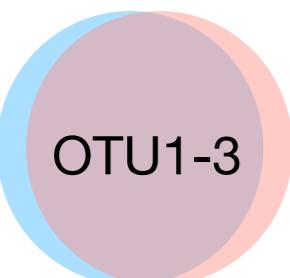
Theory



Beta-Diversity metric

Example of Bray Curtis

Same idea but considering reads, not only OTUs



$J=0$

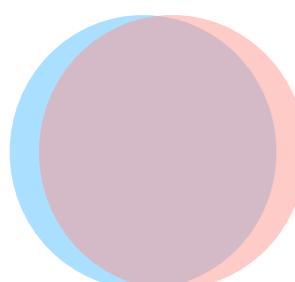
“Classical analysis” of the microbiota composition

Theory

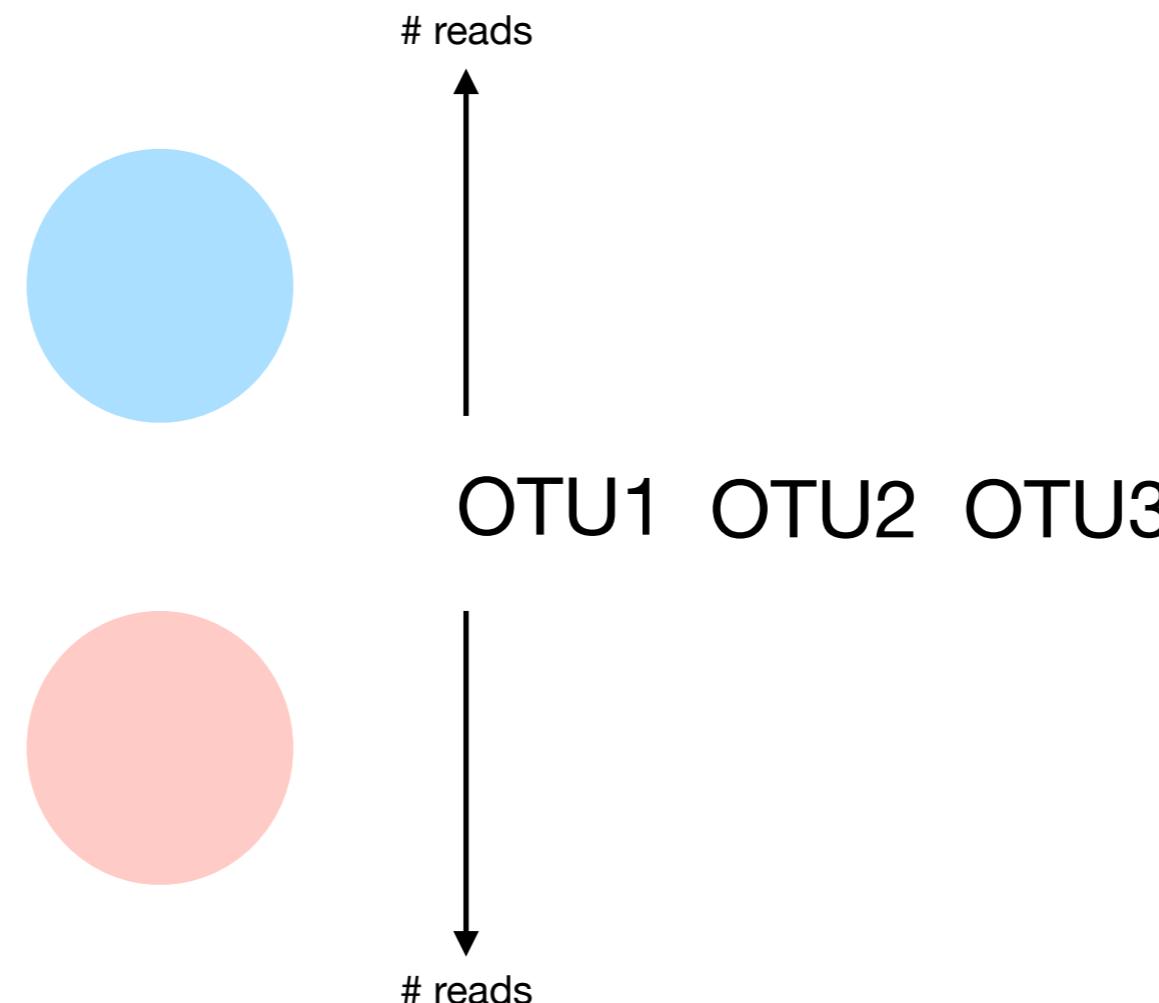


Beta-Diversity metric

Example of Bray Curtis



Same idea but considering reads, not only OTUs



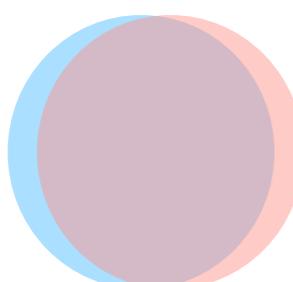
“Classical analysis” of the microbiota composition

Theory



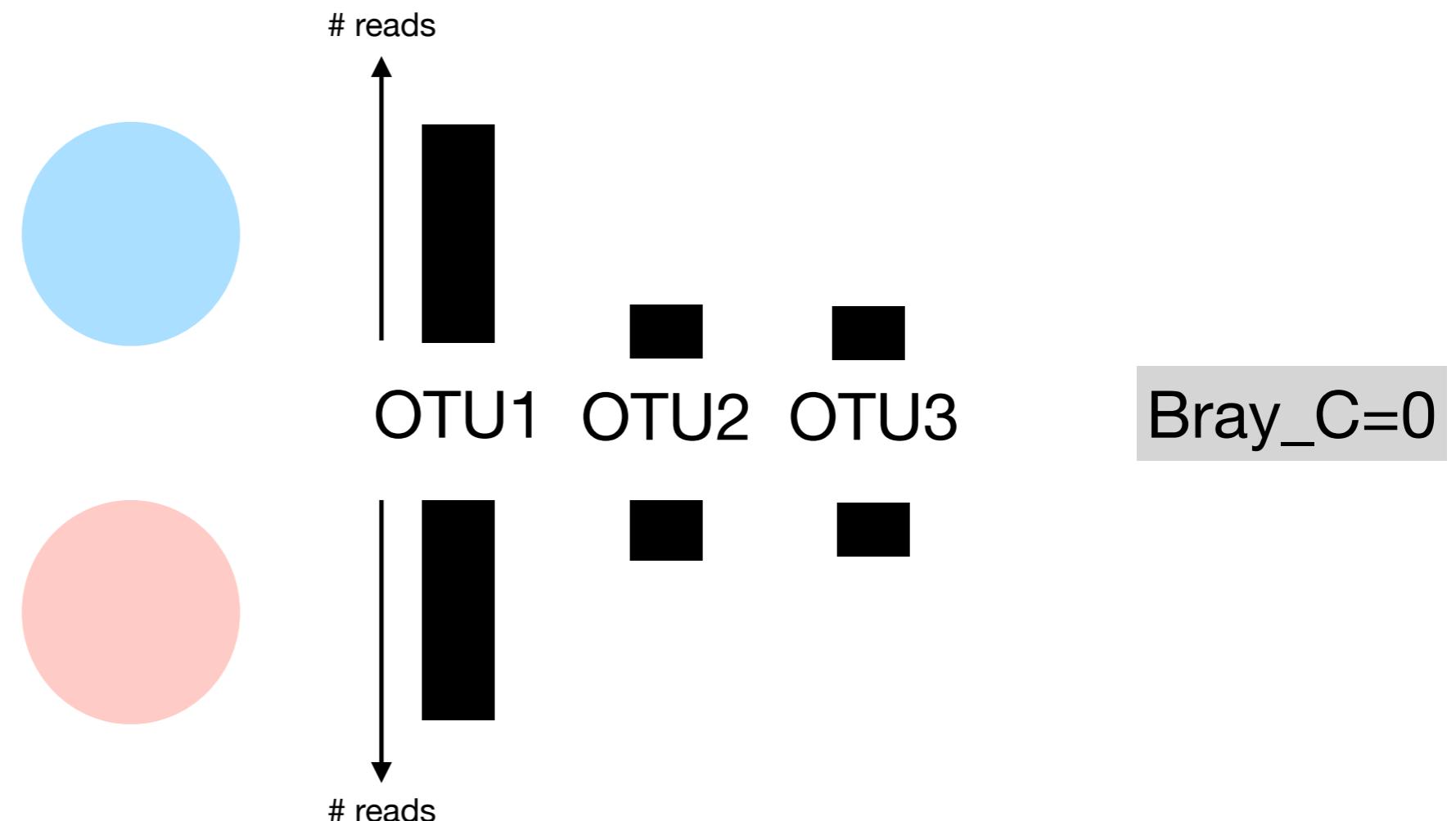
Beta-Diversity metric

Example of Bray Curtis



$J=0$

Same idea but considering reads, not only OTUs



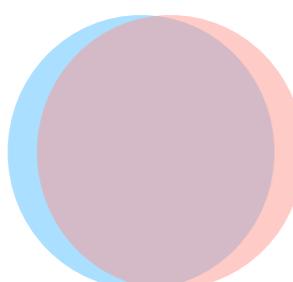
“Classical analysis” of the microbiota composition

Theory



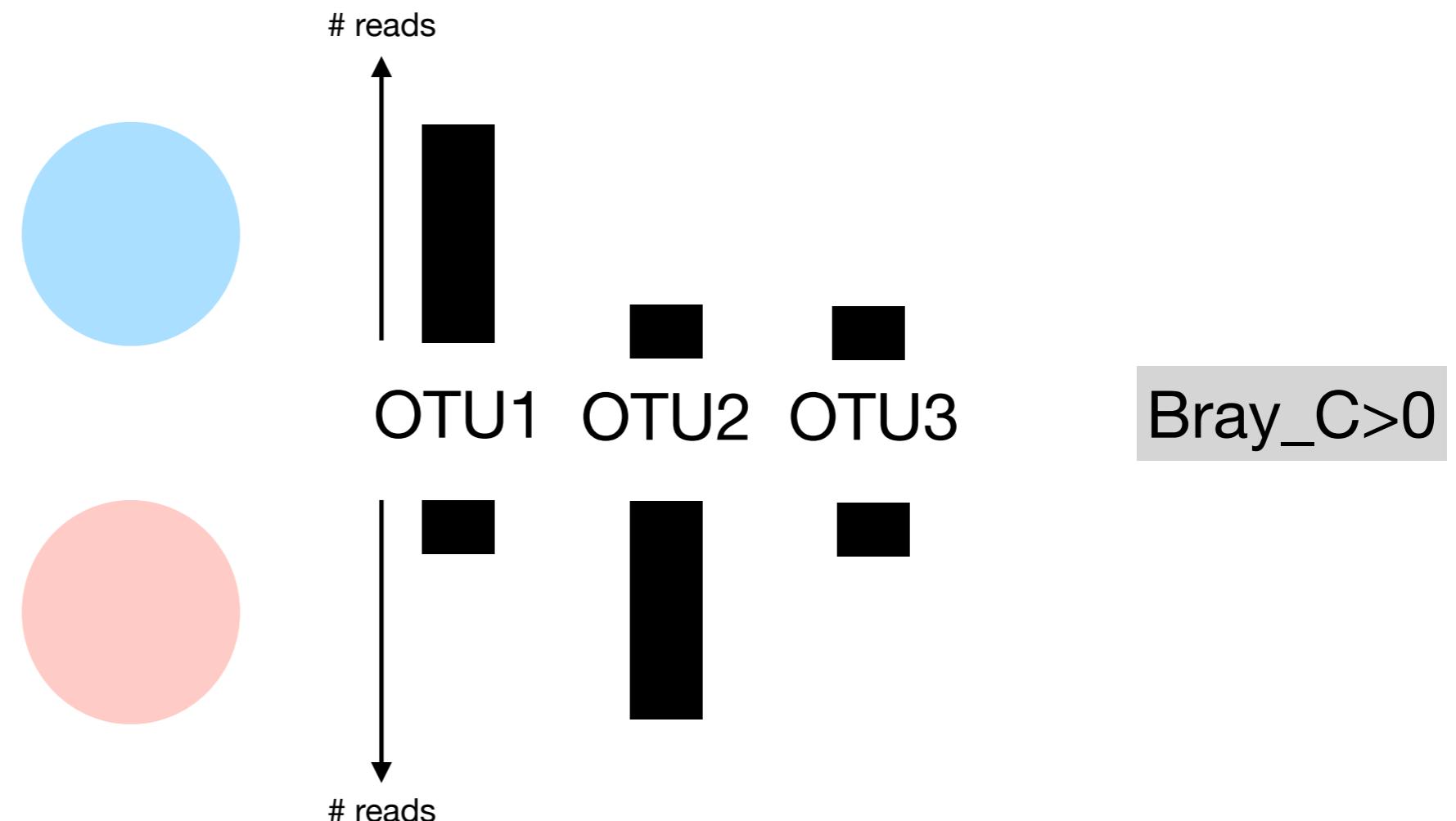
Beta-Diversity metric

Example of Bray Curtis



$J=1$

Same idea but considering reads, not only OTUs



“Classical analysis” of the microbiota composition

Theory



Beta-Diversity metric

Non – Phylogenetic

Presence/Absence

**Jaccard
Sorensen**

Abundances

BrayCurtis

“Classical analysis” of the microbiota composition

Theory



Beta-Diversity metric

Presence/Absence

Non – Phylogenetic

Phylogenetic

Jaccard
Sorensen

Abundances

BrayCurtis

“Classical analysis” of the microbiota composition

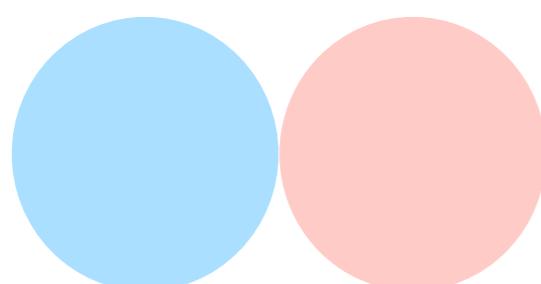
Theory



Beta-Diversity metric

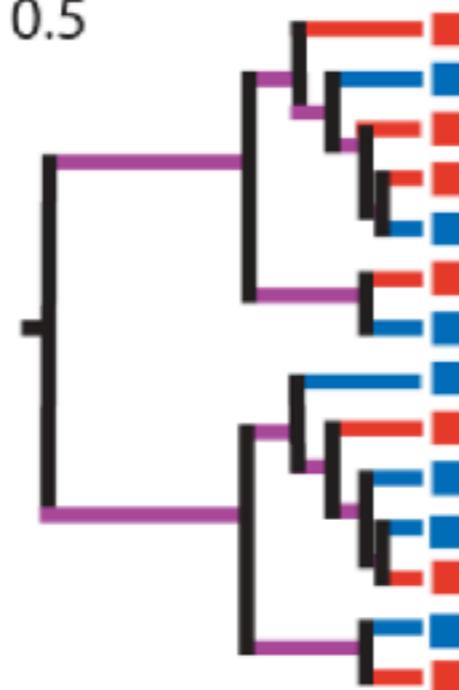
With Jaccard
(Non phylogenetic)

Example of Unifrac
(Phylogenetic-based metric)

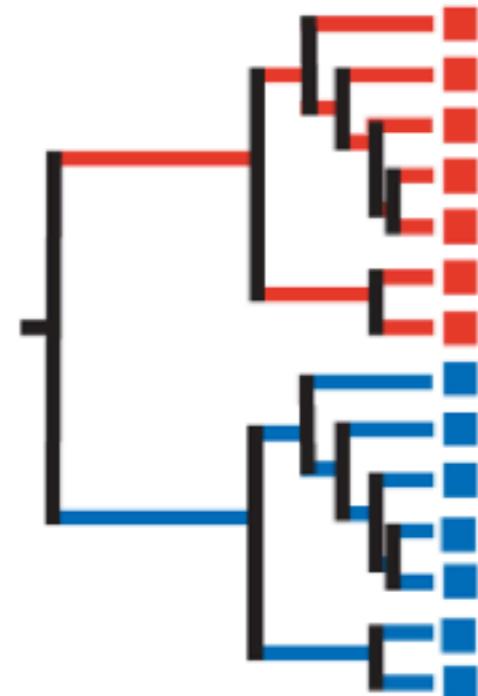


$J=1$

Related communities:
 $D \sim 0.5$

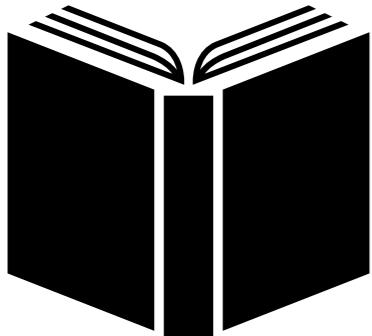


Unrelated communities:
 $D = 1$



“Classical analysis” of the microbiota composition

Theory



Beta-Diversity metric

Presence/Absence

Non – Phylogenetic

Jaccard
Sorensen

Phylogenetic

Unifrac
Phylosor

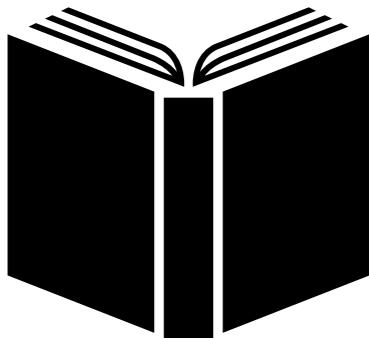
Abundances

BrayCurtis

Weighted Unifrac

“Classical analysis” of the microbiota composition

Theory



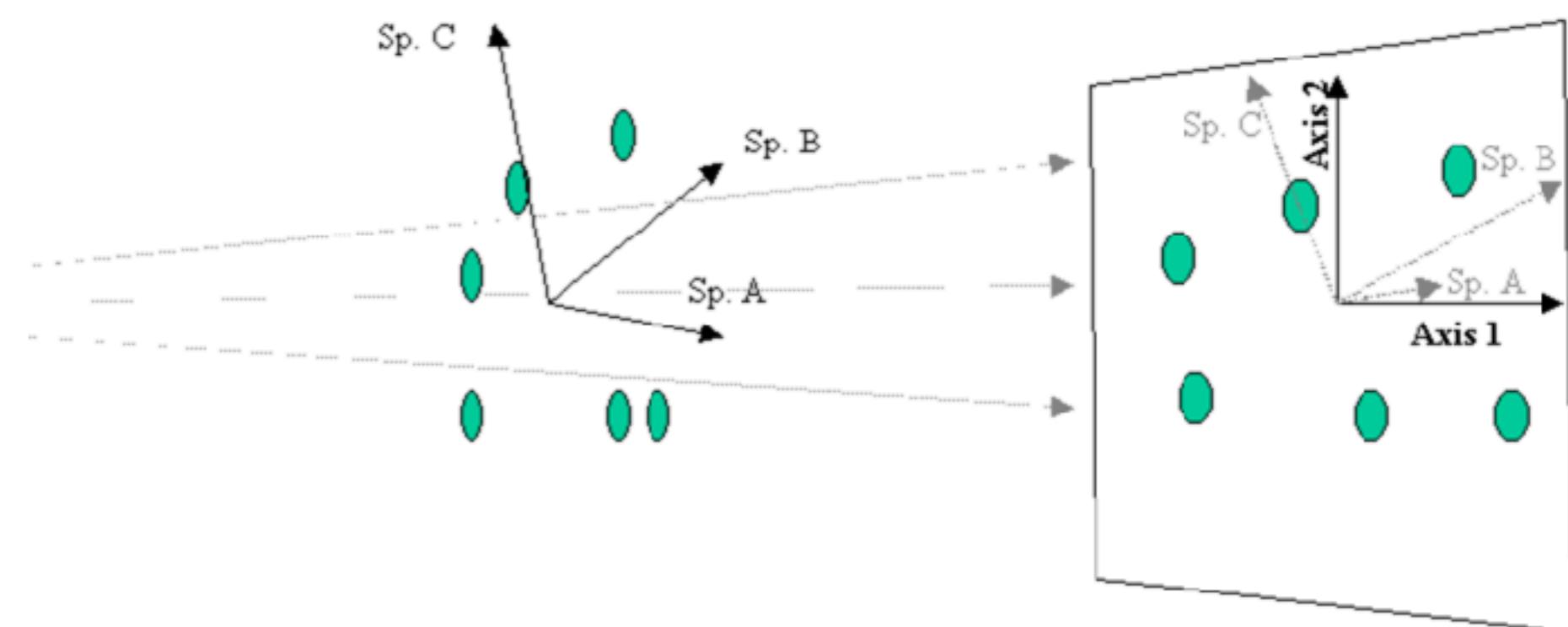
Visualization

Try to best represent these distances on a limited number of axes

for n (number of OTUs) dimensions to 2 (or 3) dimensions

Three
species

Two
dimensions

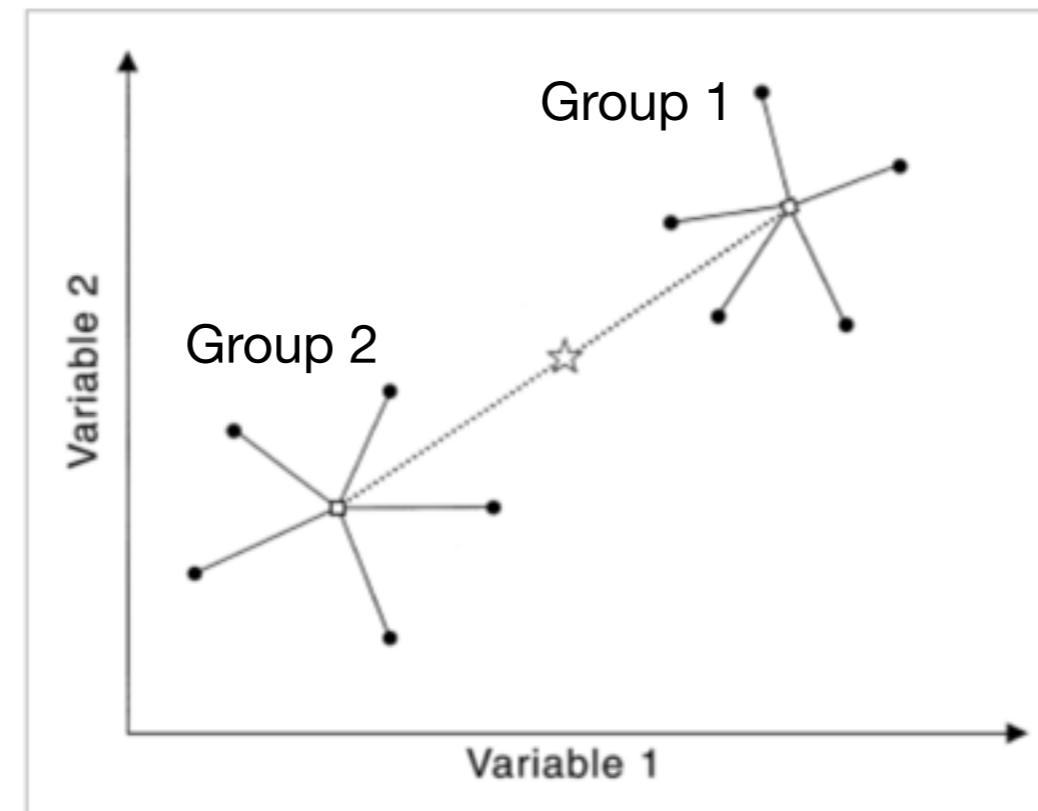


“Classical analysis” of the microbiota composition

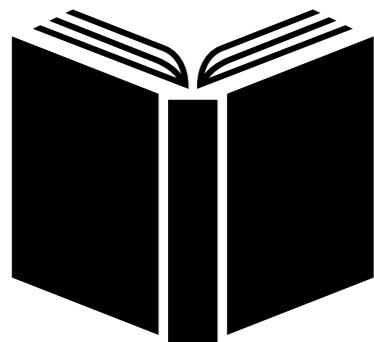


Statistical test

Example of the PERMANOVA test

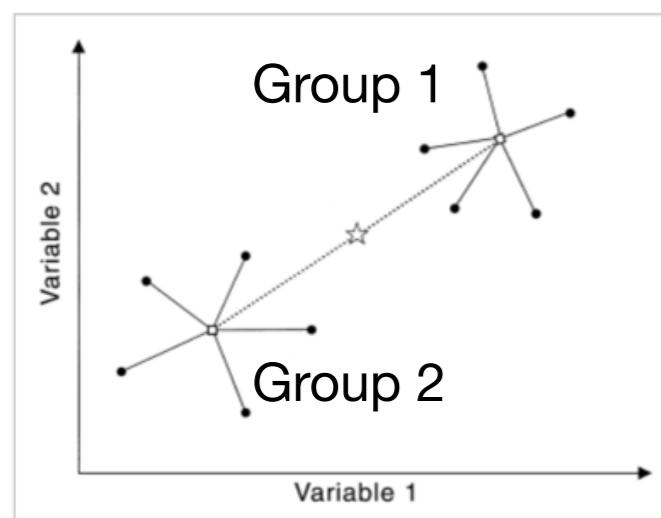


“Classical analysis” of the microbiota composition



Statistical test

Example of the PERMANOVA test



Beta-div

Beta-div

Observations

(symmetric)

DISTANCE MATRIX

Observations

Group

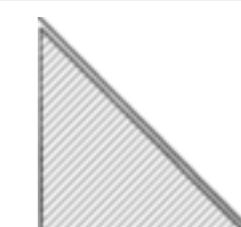
Group 2

Group

Group

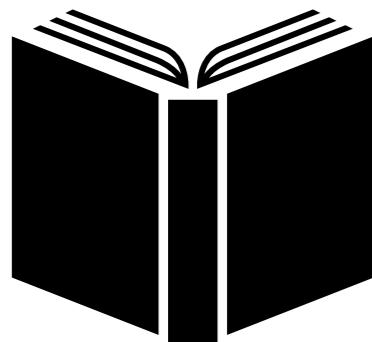
[between]

F ratio = _____



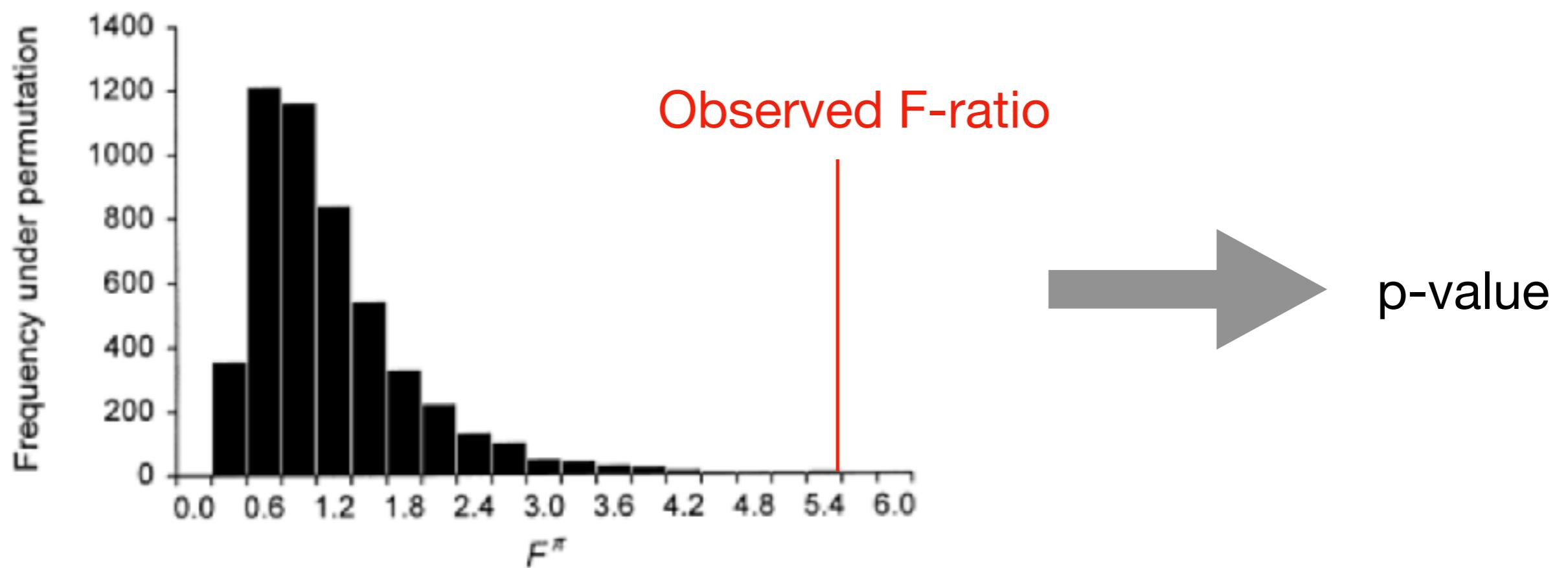
[within]

“Classical analysis” of the microbiota composition



Statistical test

Example of the PERMANOVA test



“Classical analysis” of the microbiota composition



Statistical test

Example of the PERMANOVA test

Call:

```
adonis(formula = UniFracBeta ~ Depth_m, data = data.frame(sample_data(mothur)))
```

Permutation: free

Number of permutations: 999

Terms added sequentially (first to last)

	Df	SumsOfSqs	MeanSqs	F.Model	R2	Pr(>F)
Depth_m	1	0.34826	0.34826	2.5771	0.34012	0.001 ***
Residuals	5	0.67569	0.13514		0.65988	
Total	6	1.02395			1.00000	

Signif. codes: 0 ‘***’ 0.001 ‘**’ 0.01 ‘*’ 0.05 ‘.’ 0.1 ‘ ’ 1

3. Analysis the phylogenetic composition of microbiomes

Practice



Code lines 295 – 391

Import data

```
> read.table()  
> read.tree()
```

Create a unified object

```
> phyloseq()
```

Compute “traditional”
Beta-diversity metrics

```
> vegdist()  
> UniFrac()
```

Plot their relationships

```
> plot()
```

Visualize Beta-diversities

```
> ordinate()
```

Statistical test

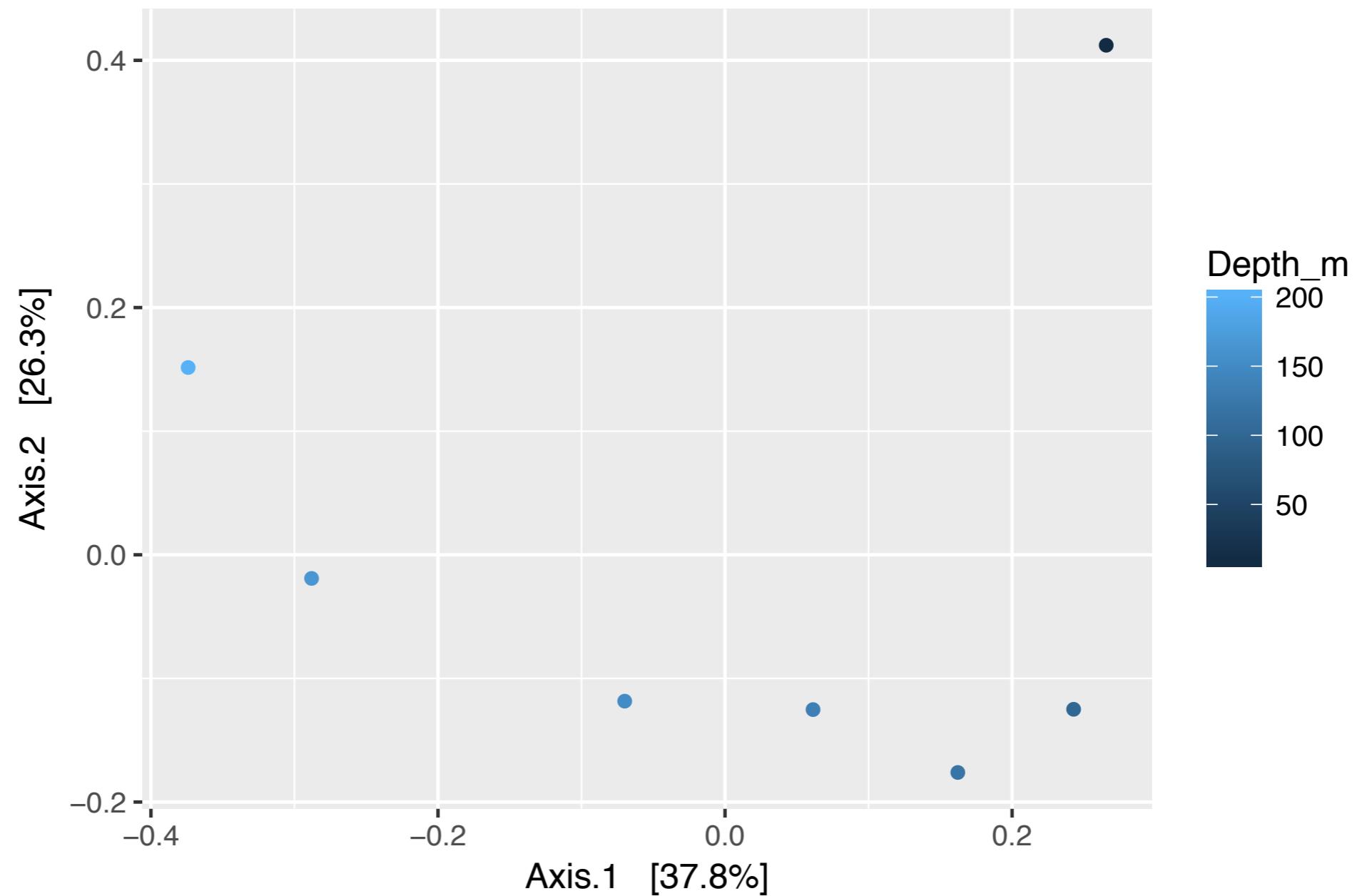
```
> adonis()
```

“Classical analysis” of the microbiota composition

Practice



Visualization



“Classical analysis” of the microbiota composition

Practice



Statistical test

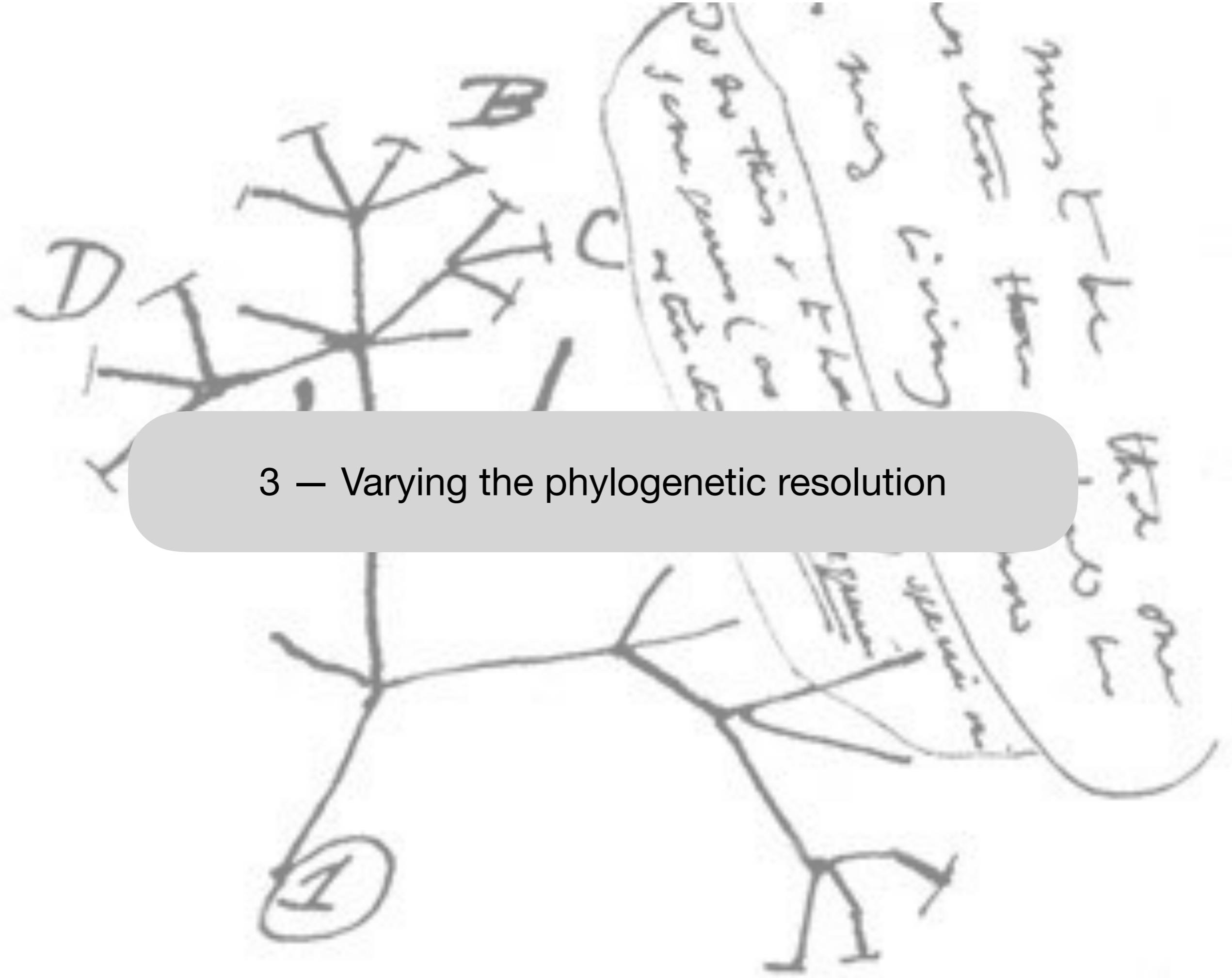
R function

```
> adonis()
```

Test several beta-diversity metrics: difference between phylogenetic / non phylogenetic?

Test one or more explanatory variables

3 – Varying the phylogenetic resolution



3 – Varying the phylogenetic resolution

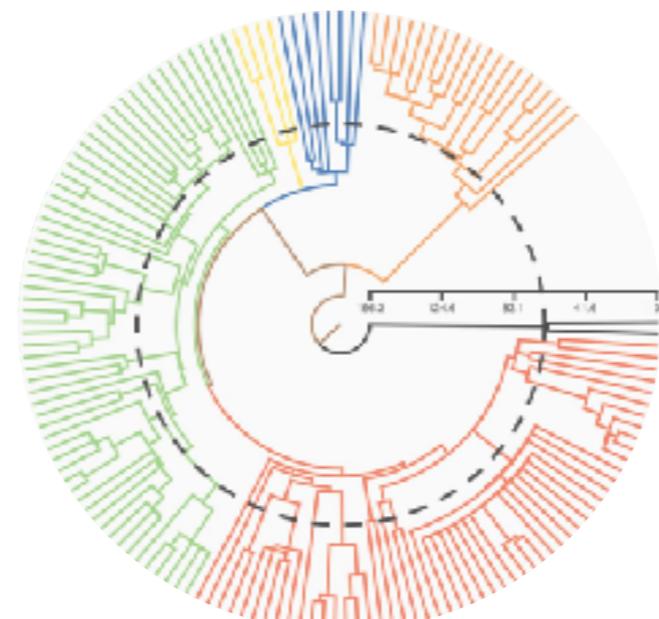
Example 1: Mammalian guts

GUT MICROBIOTA
COMPOSITION

~ HOST PHYLOGENY

OR

HOST DIET ?



Microb Ecol
DOI 10.1007/s00248-017-1041-8

HOST MICROBE INTERACTIONS



CrossMark

**Diet Versus Phylogeny: a Comparison of Gut
Microbiota in Captive Colobine Monkey Species**

3 – Varying the phylogenetic resolution

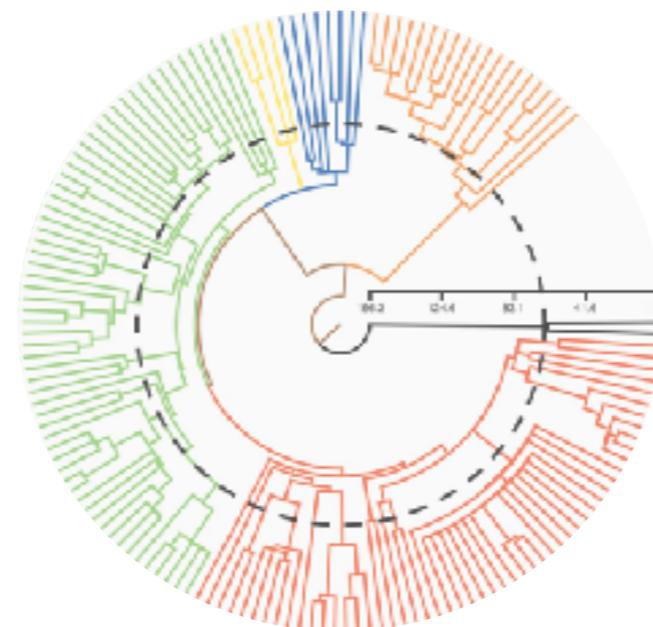
Example 1: Mammalian guts

GUT MICROBIOTA
COMPOSITION

~ HOST PHYLOGENY

OR

HOST DIET ?

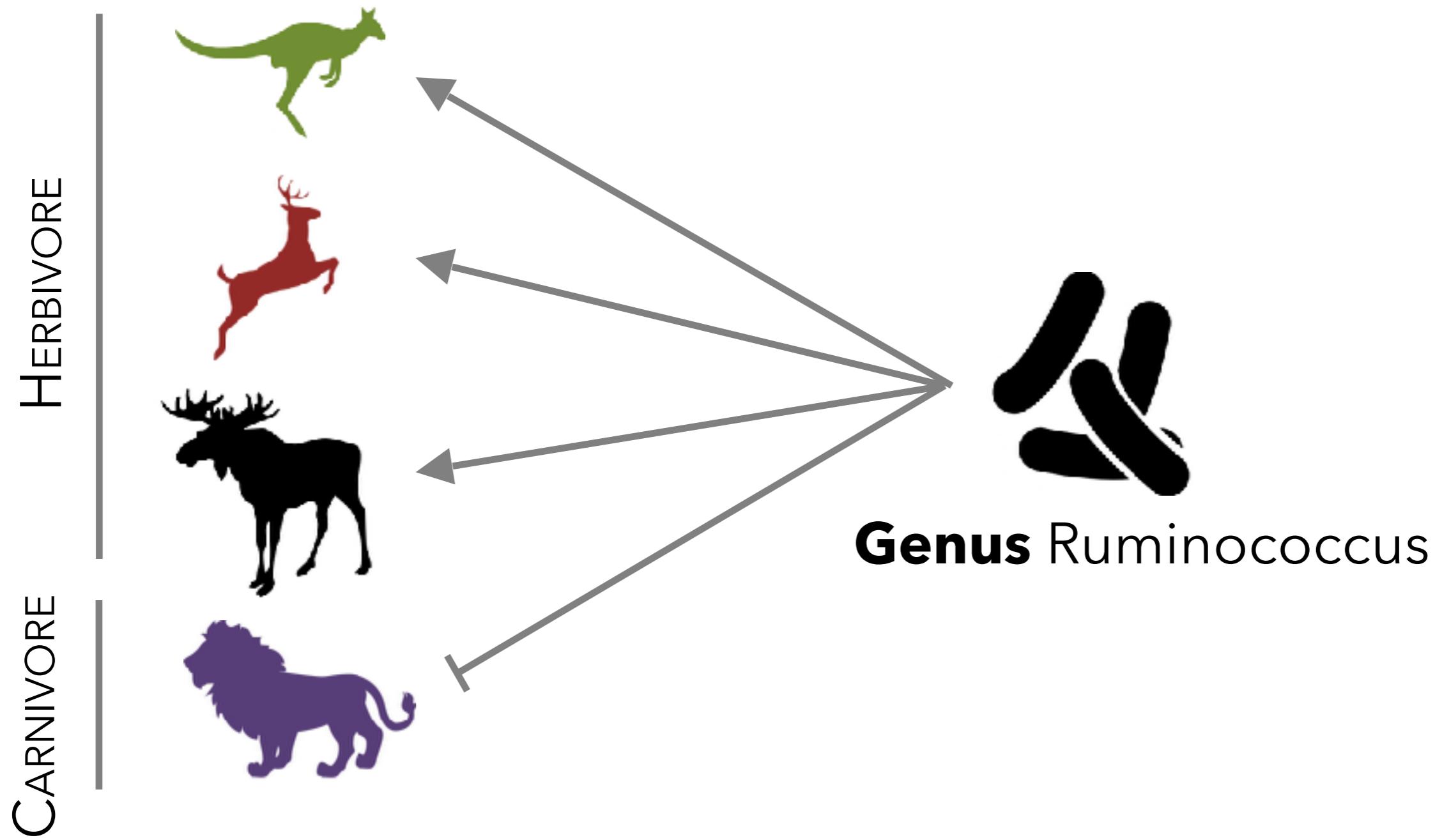


Idea: Different factors might play at different phylogenetic resolutions

3 – Varying the phylogenetic resolution

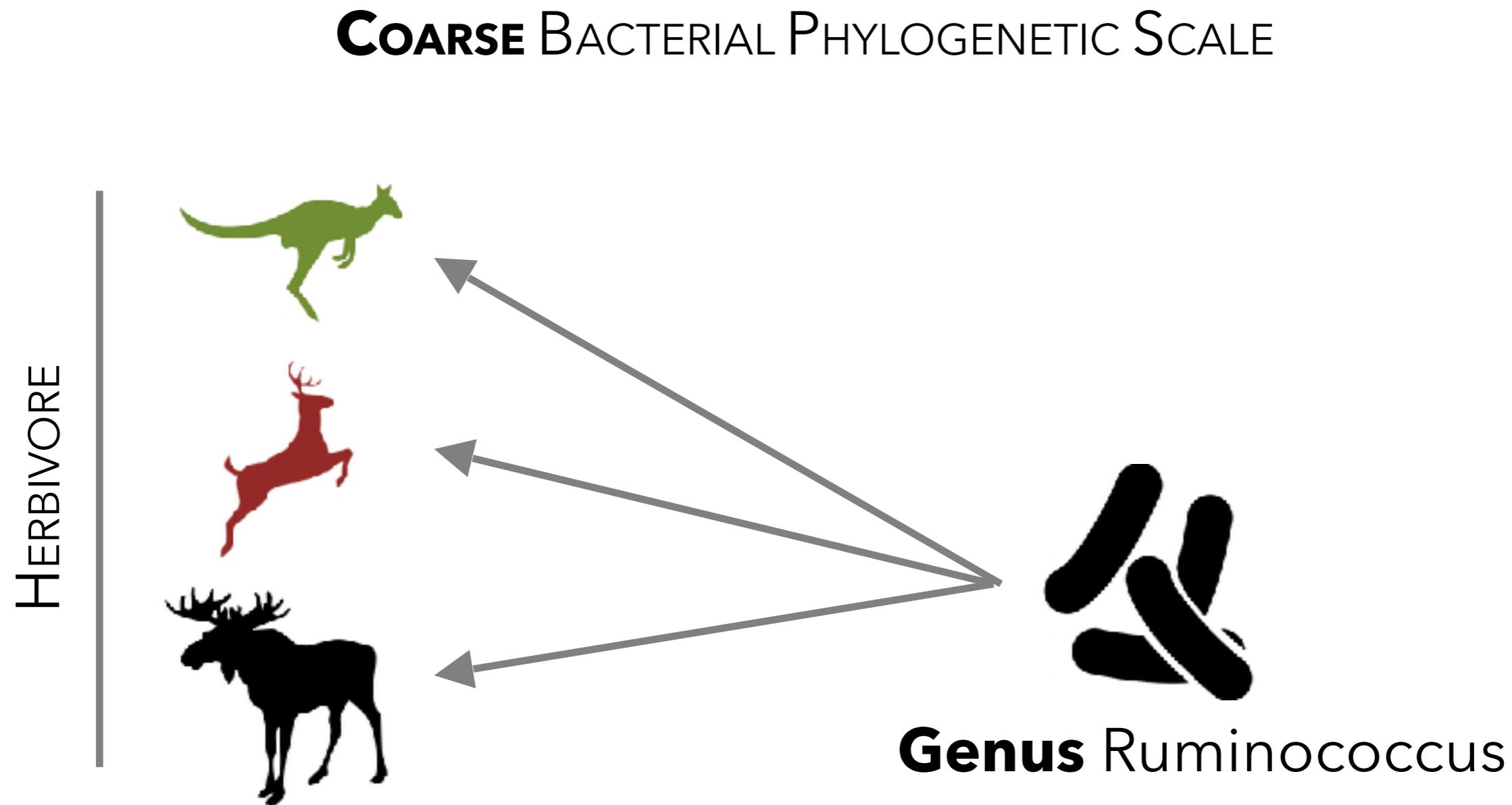
Example 1: Mammalian guts

COARSE BACTERIAL PHYLOGENETIC SCALE



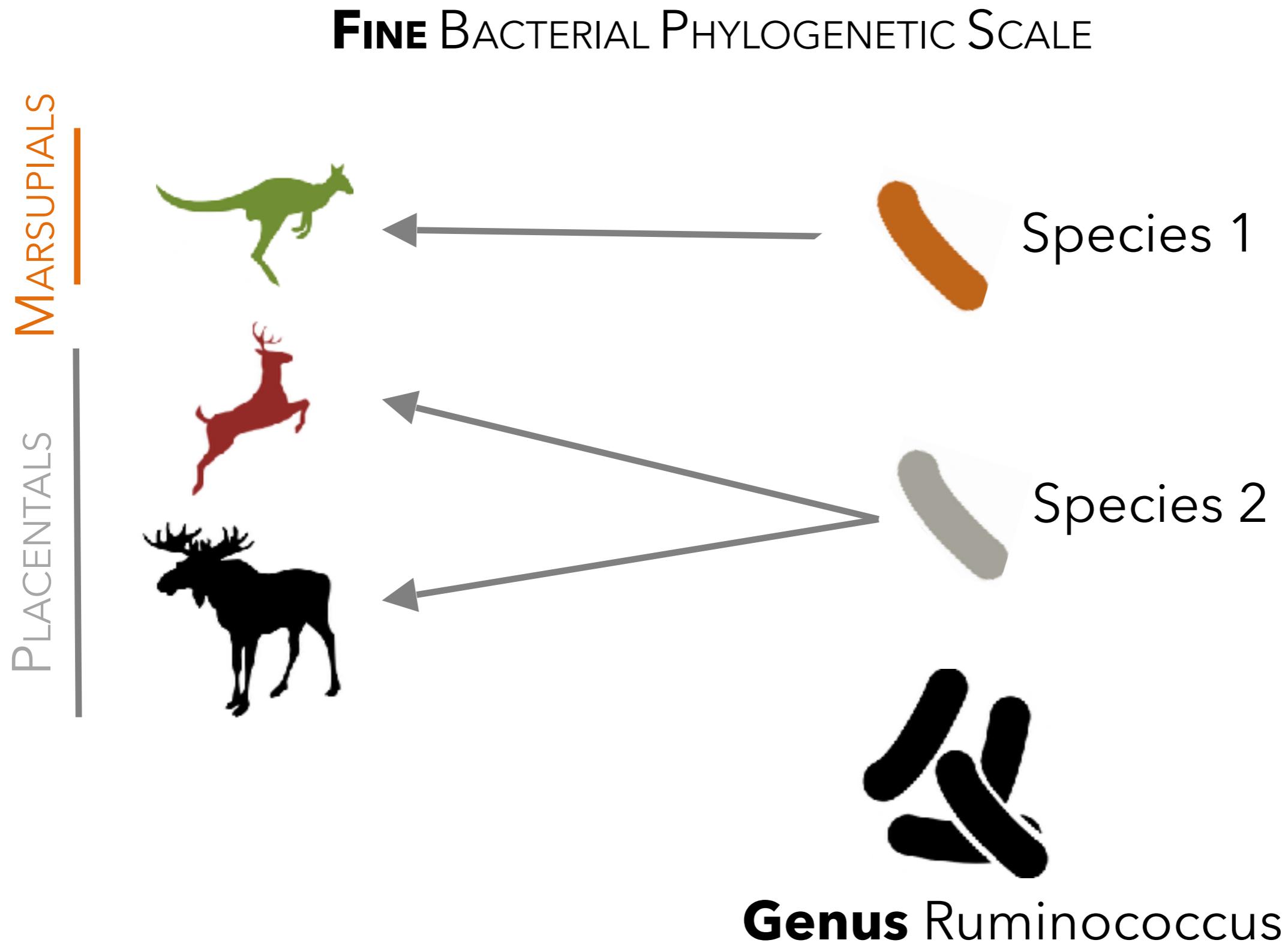
3 – Varying the phylogenetic resolution

Example 1: Mammalian guts

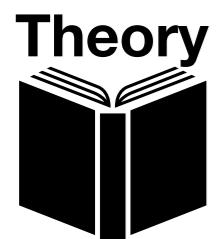


3 – Varying the phylogenetic resolution

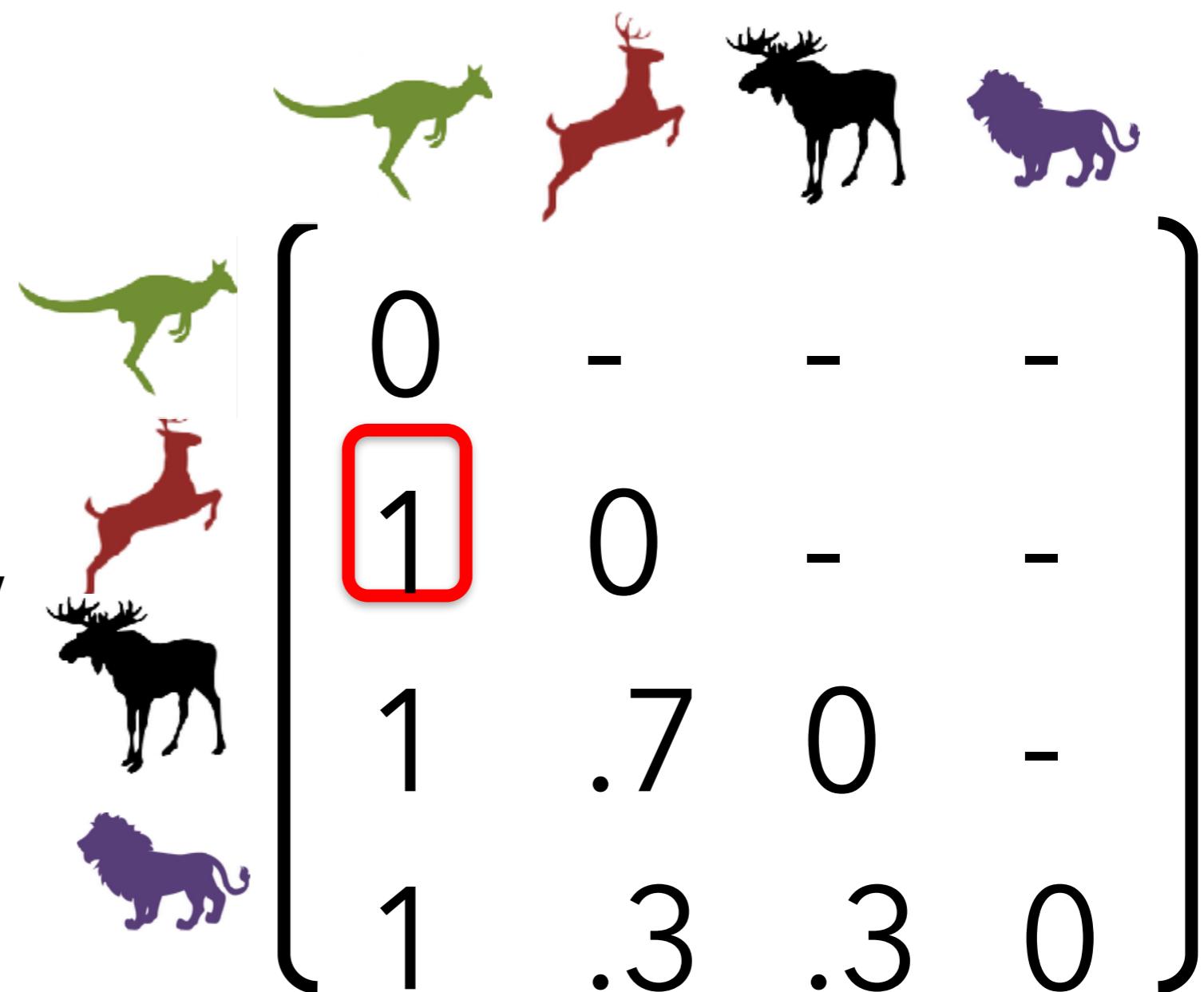
Example 1: Mammalian guts



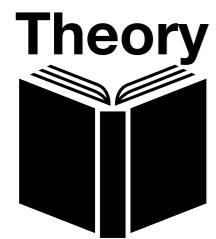
3 – Varying the phylogenetic resolution



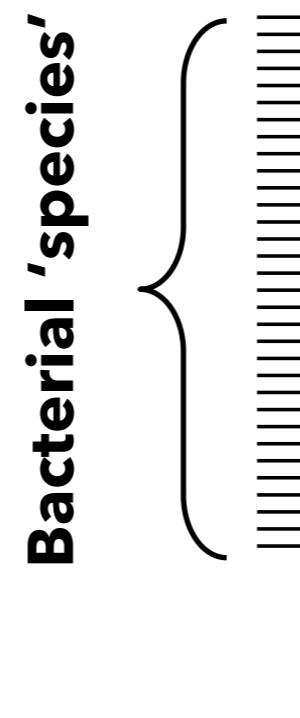
% of shared **'species'**
between hosts



3 – Varying the phylogenetic resolution



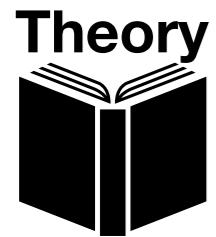
Method: “Slicing the tree of life”



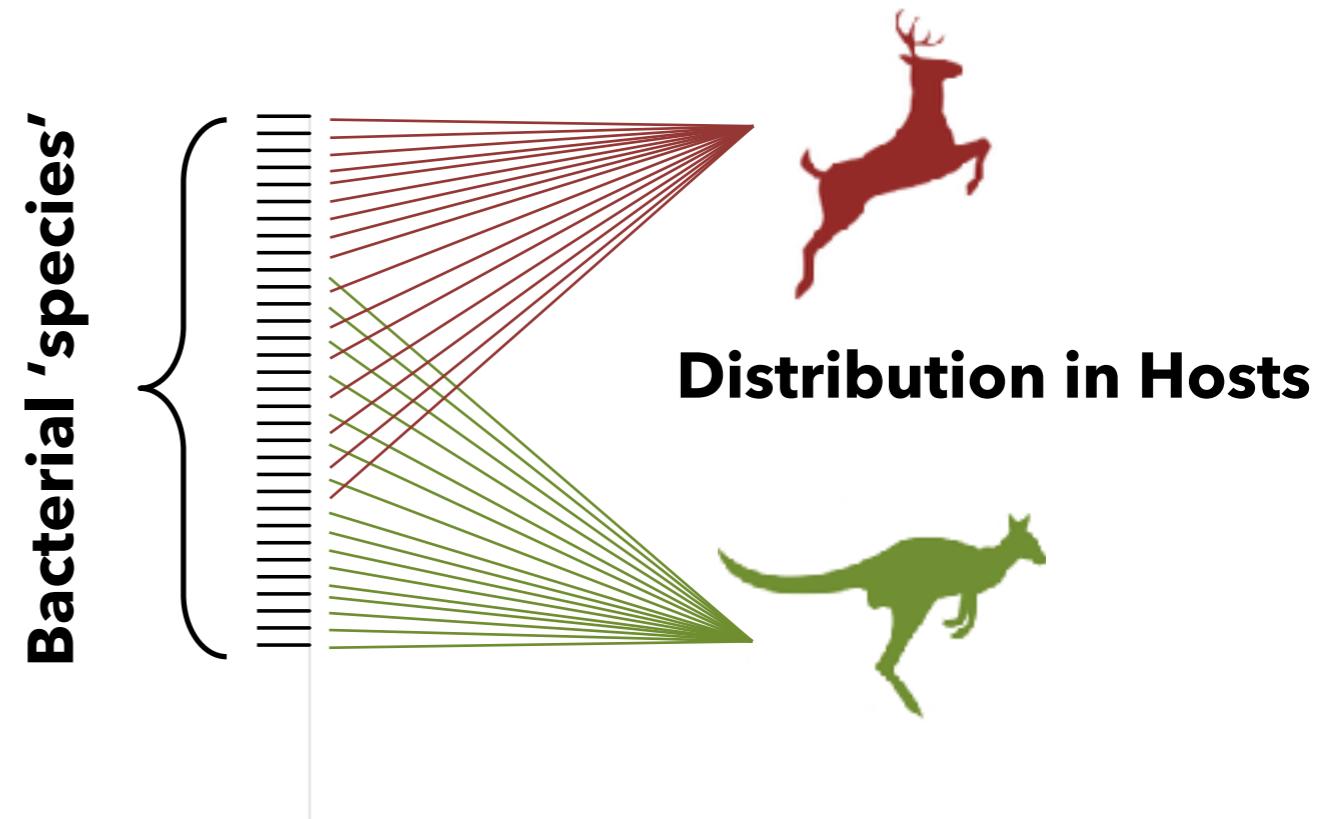
Distribution in Hosts



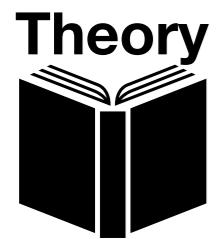
3 – Varying the phylogenetic resolution



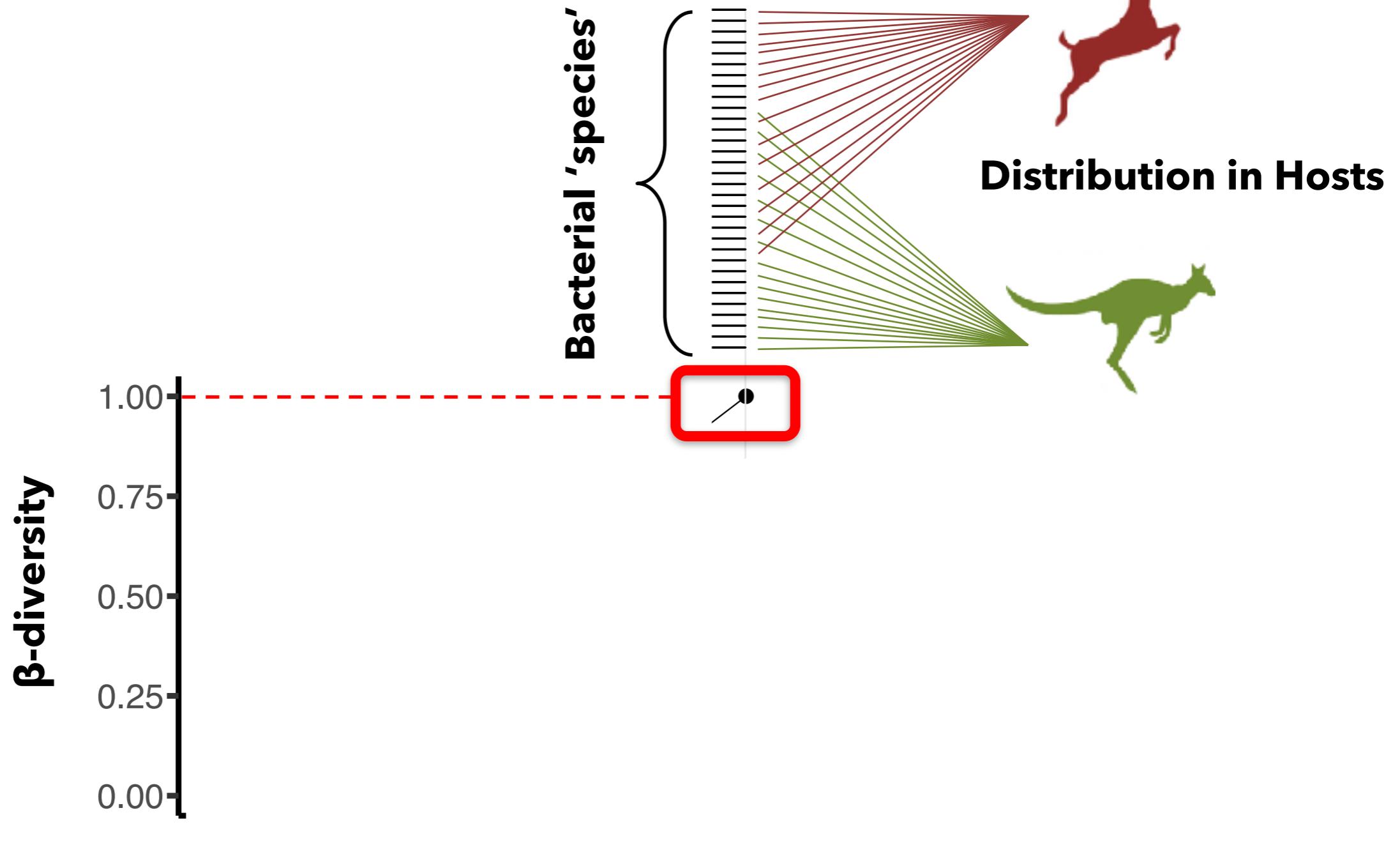
Method: “Slicing the tree of life”



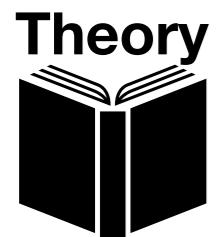
3 – Varying the phylogenetic resolution



Method: “Slicing the tree of life”

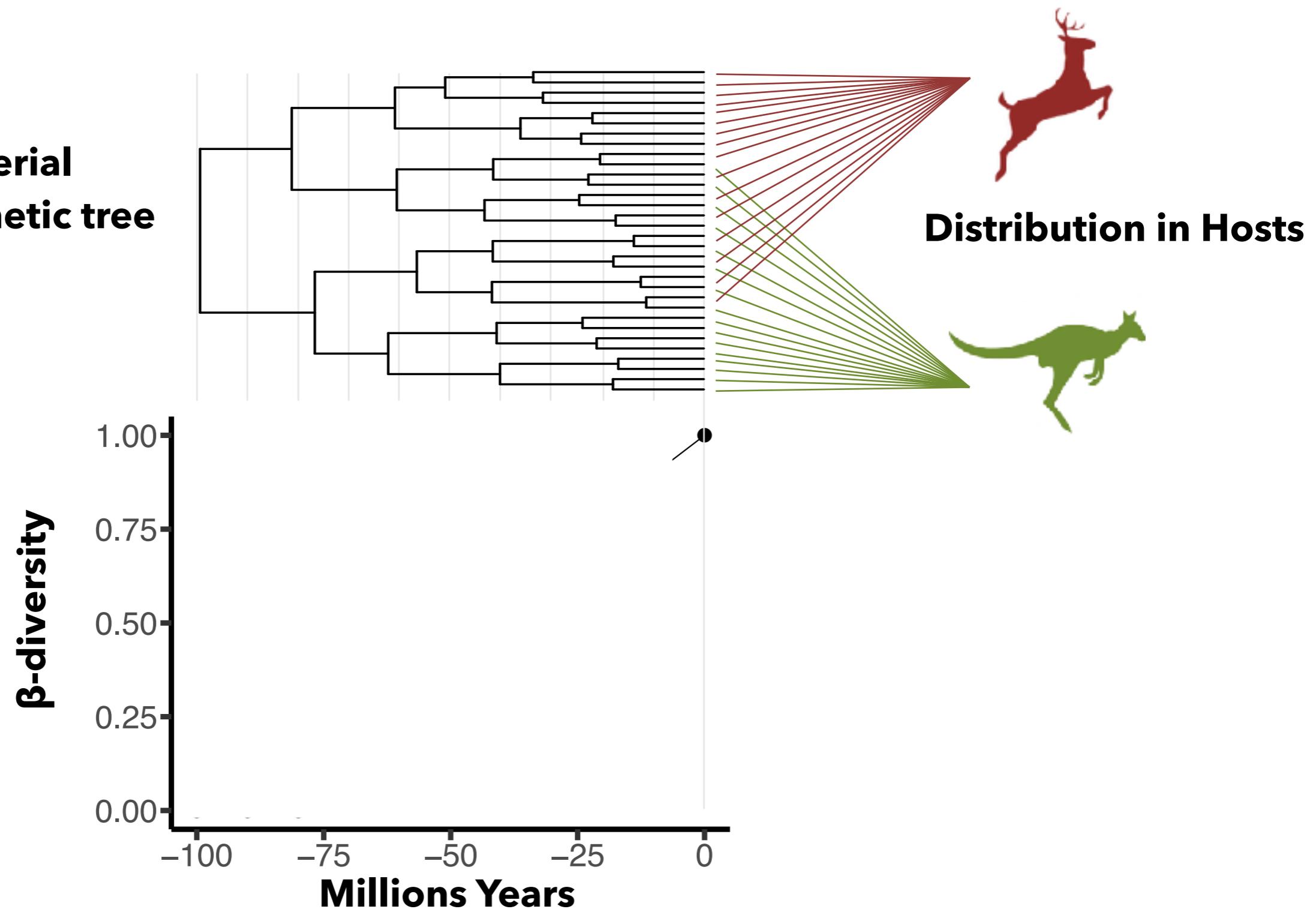


3 – Varying the phylogenetic resolution

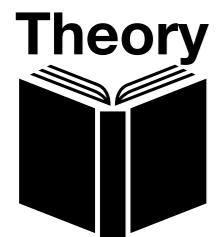


Method: “Slicing the tree of life”

**Bacterial
Phylogenetic tree**

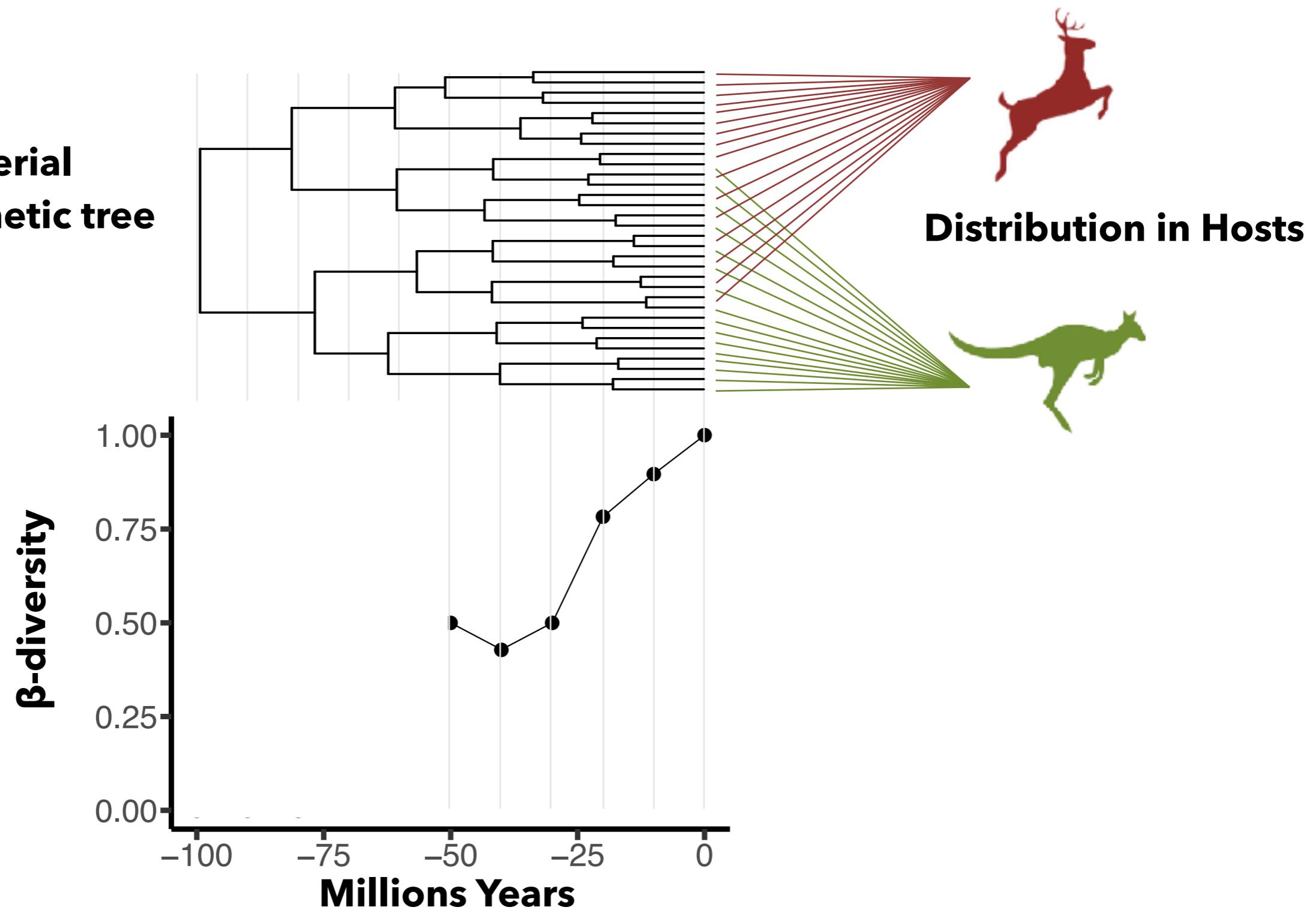


3 – Varying the phylogenetic resolution

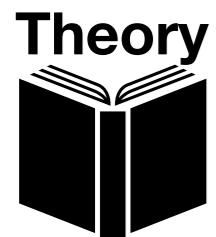


Method: “Slicing the tree of life”

**Bacterial
Phylogenetic tree**

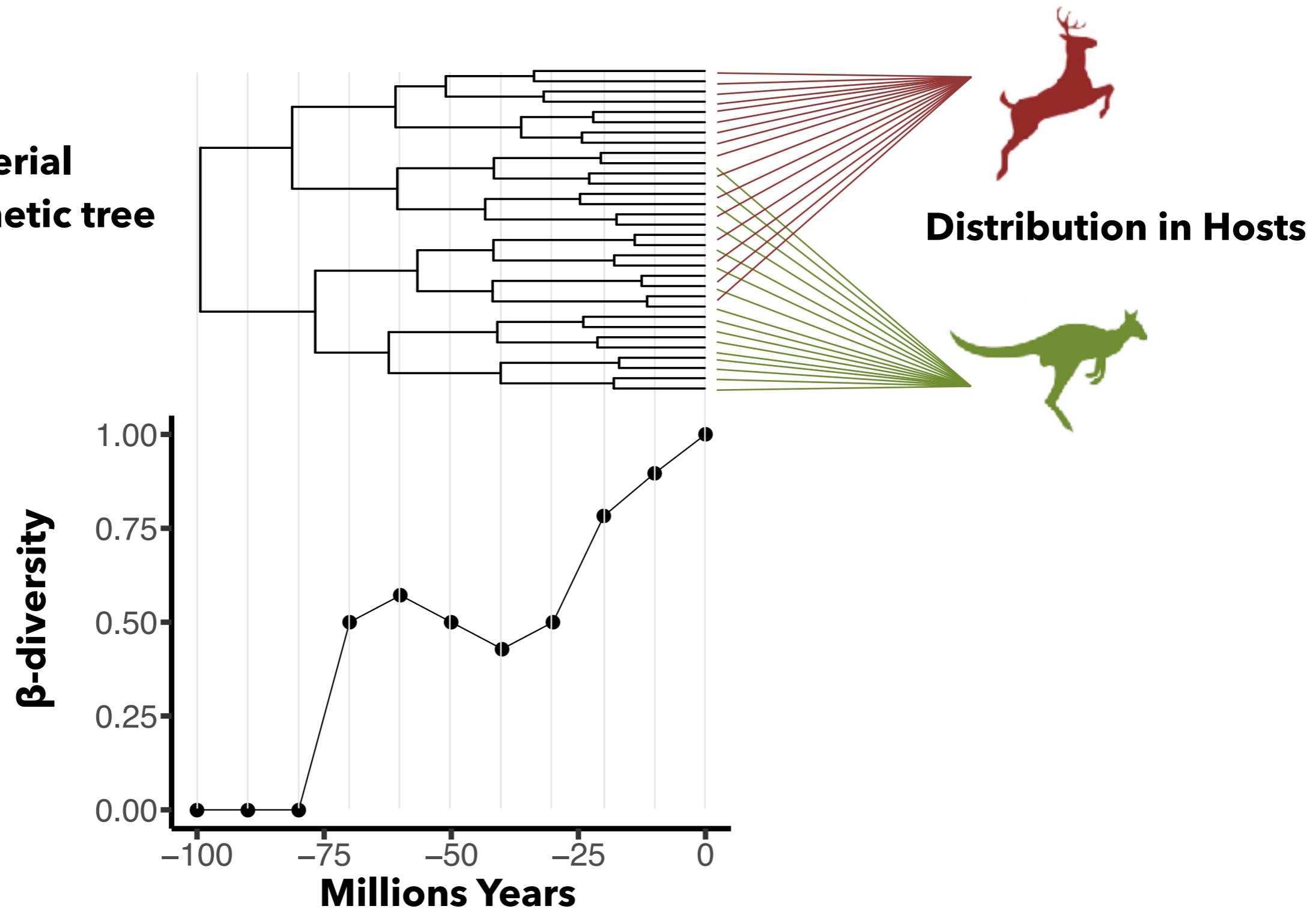


3 – Varying the phylogenetic resolution

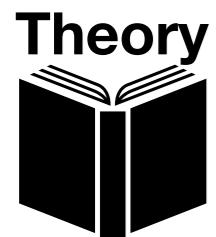


Method: “Slicing the tree of life”

**Bacterial
Phylogenetic tree**

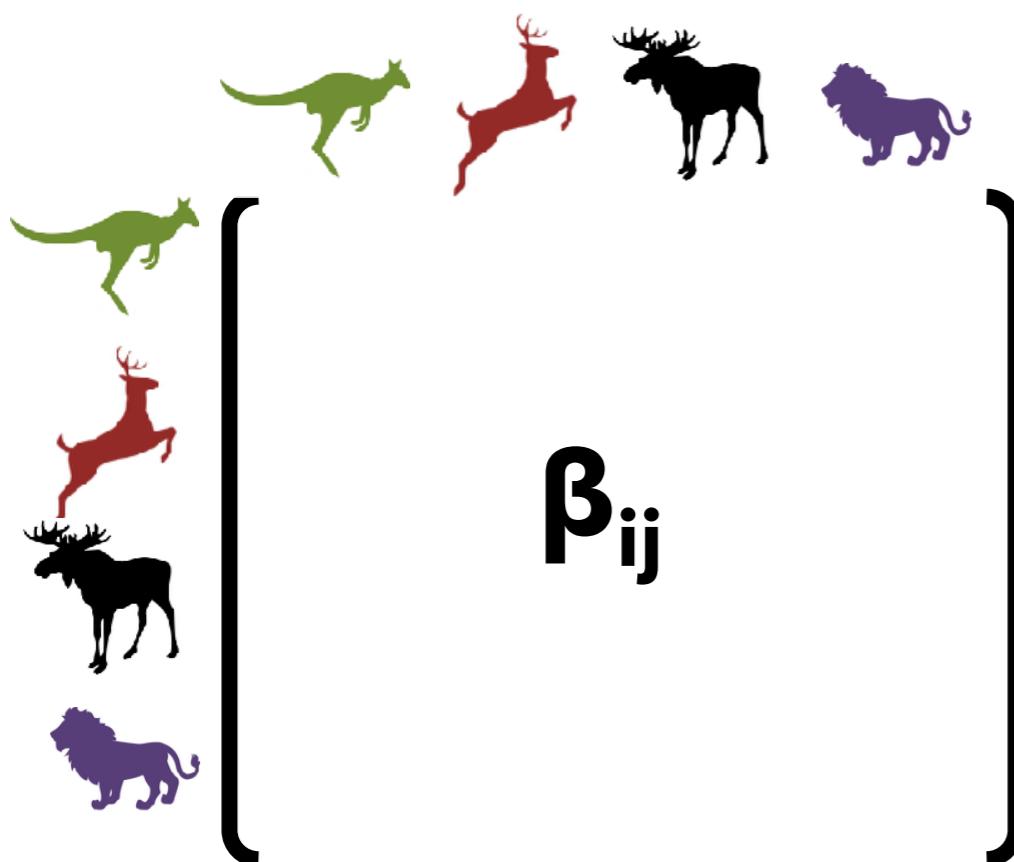


3 – Varying the phylogenetic resolution

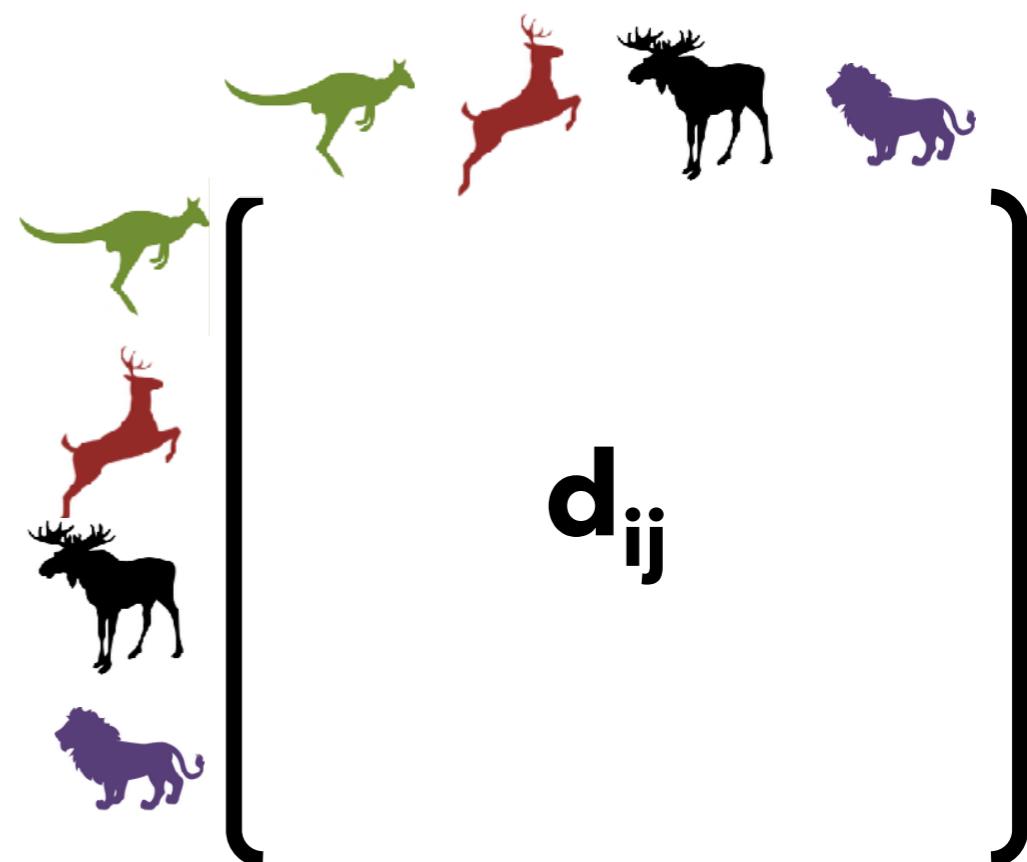


Method: “Slicing the tree of life”

β-DIVERSITY

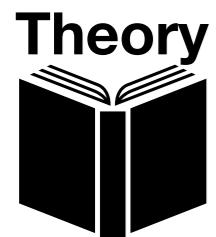


HOST PHYLO DISTANCES **HOST DIET DISTANCES**



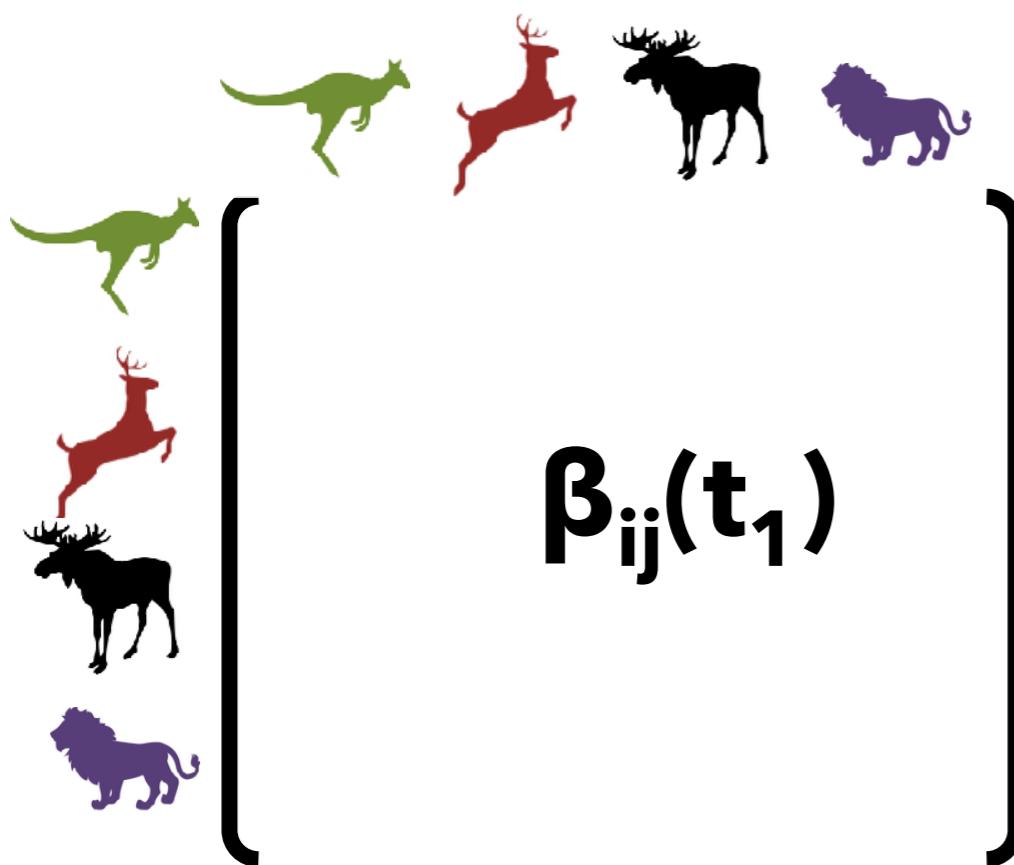
~

3 – Varying the phylogenetic resolution



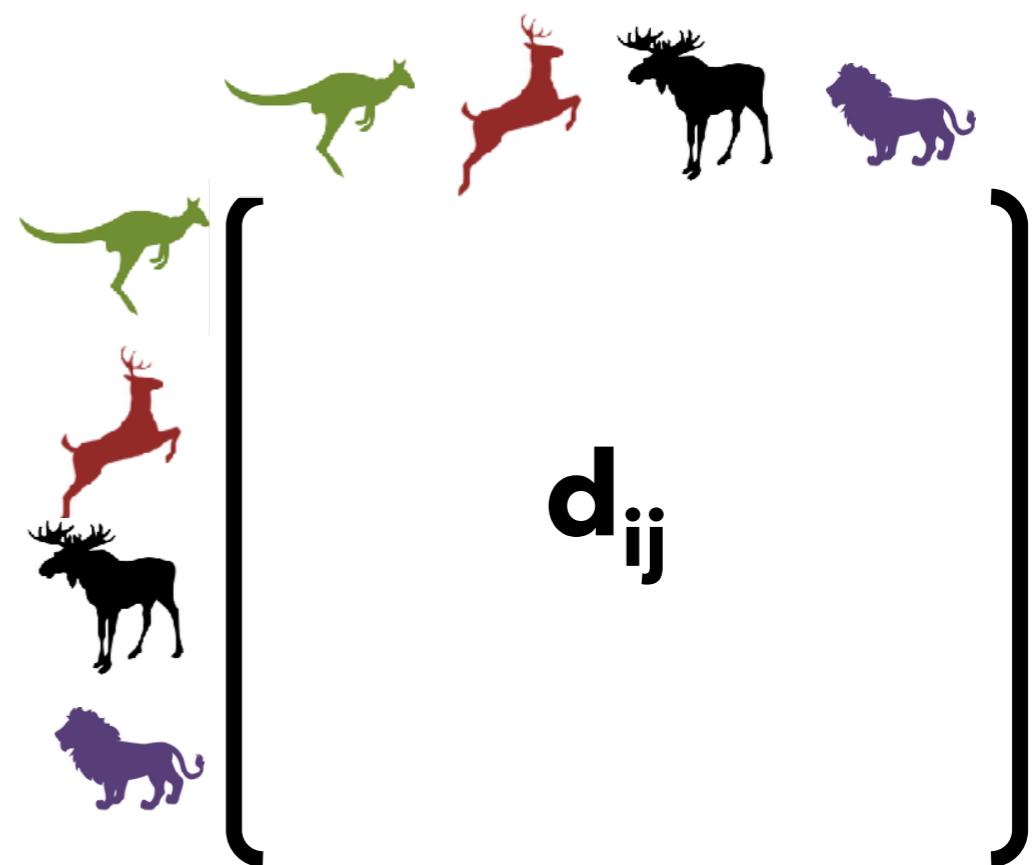
Method: “Slicing the tree of life”

β-DIVERSITY

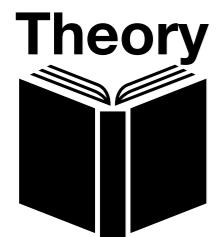


~

HOST PHYLO DISTANCES **HOST DIET DISTANCES**

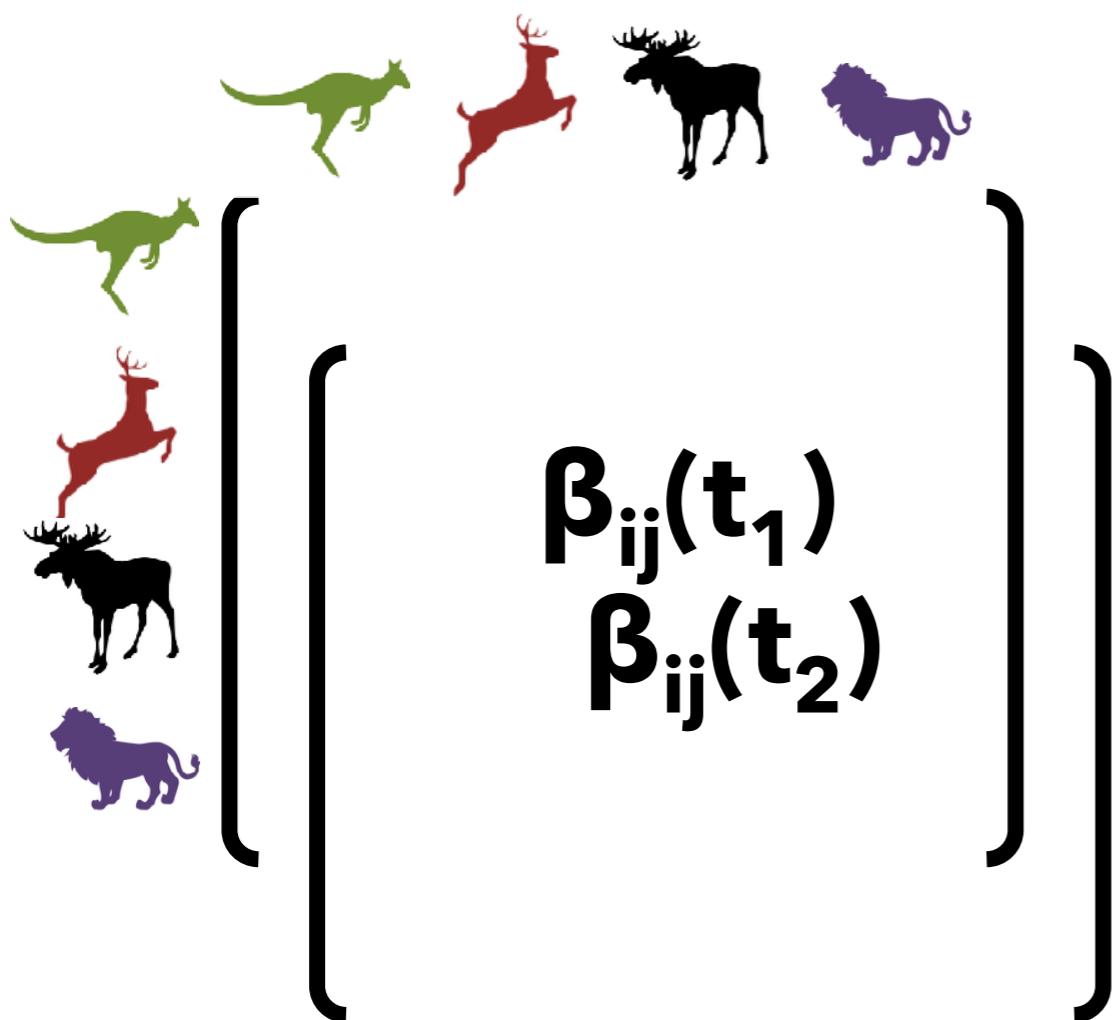


3 – Varying the phylogenetic resolution



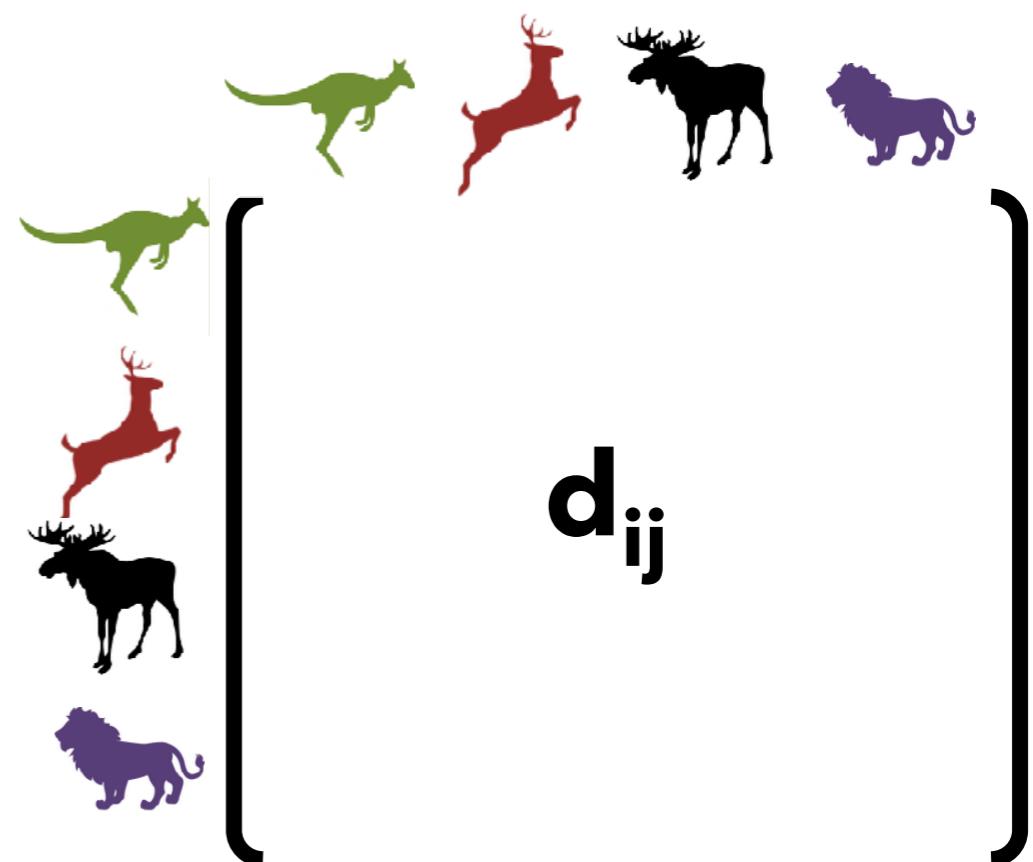
Method: “Slicing the tree of life”

β-DIVERSITY

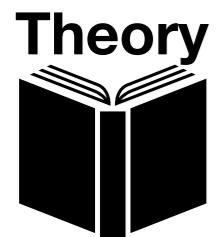


HOST PHYLO DISTANCES

HOST DIET DISTANCES



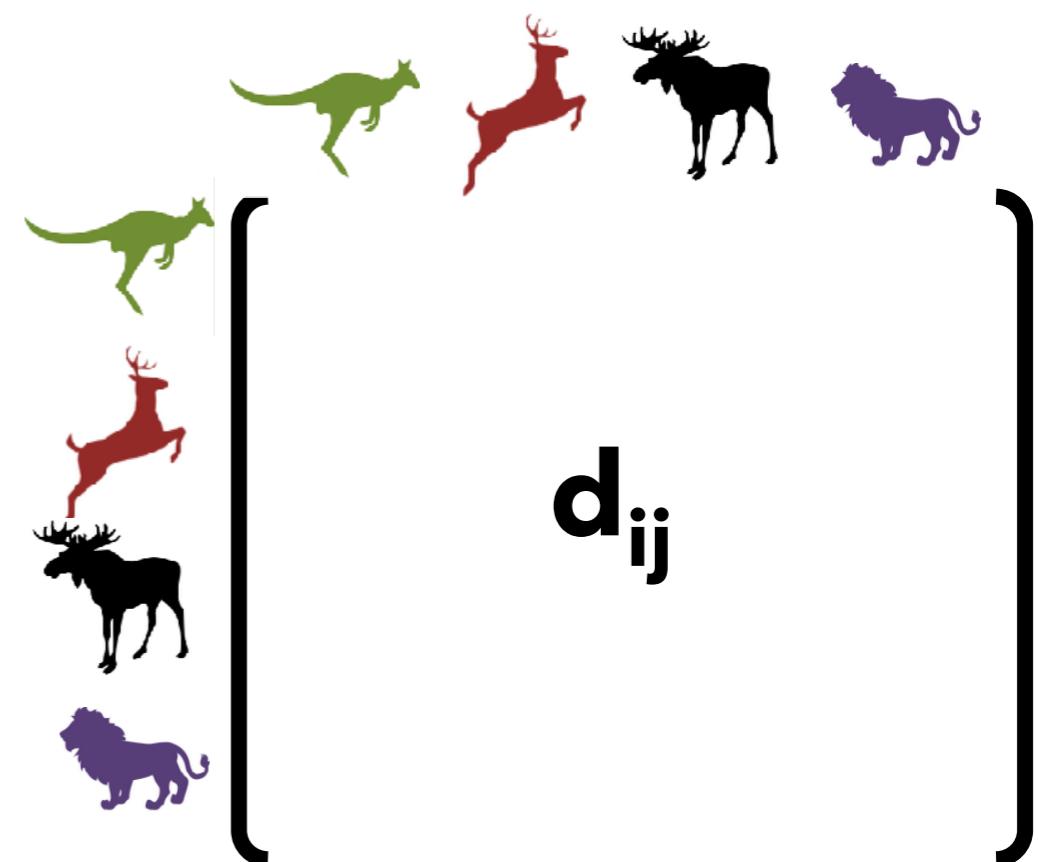
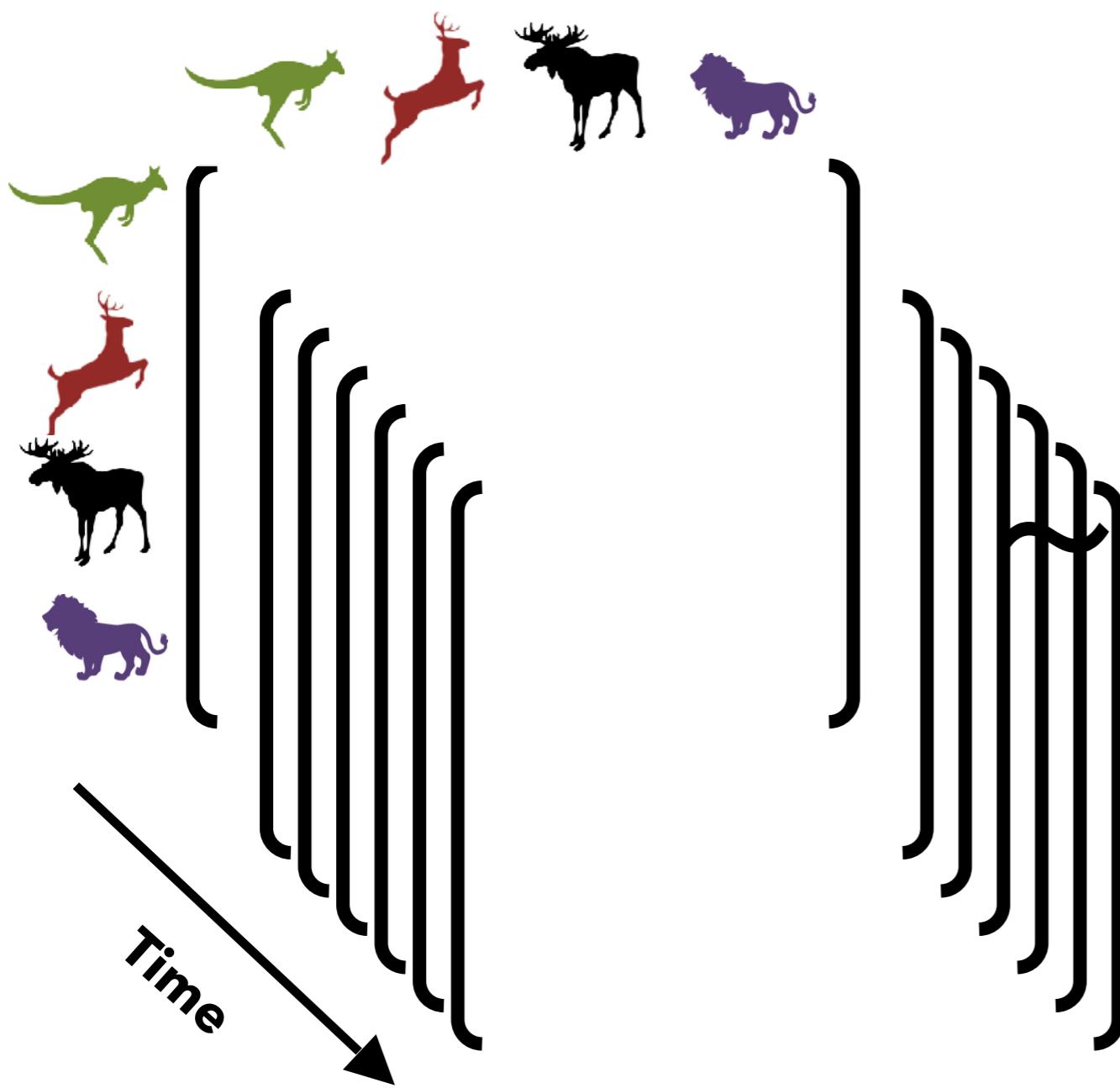
3 – Varying the phylogenetic resolution



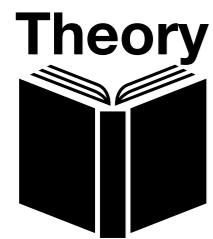
Method: “Slicing the tree of life”

HOST PHYLO DISTANCES
HOST DIET DISTANCES

β-DIVERSITY

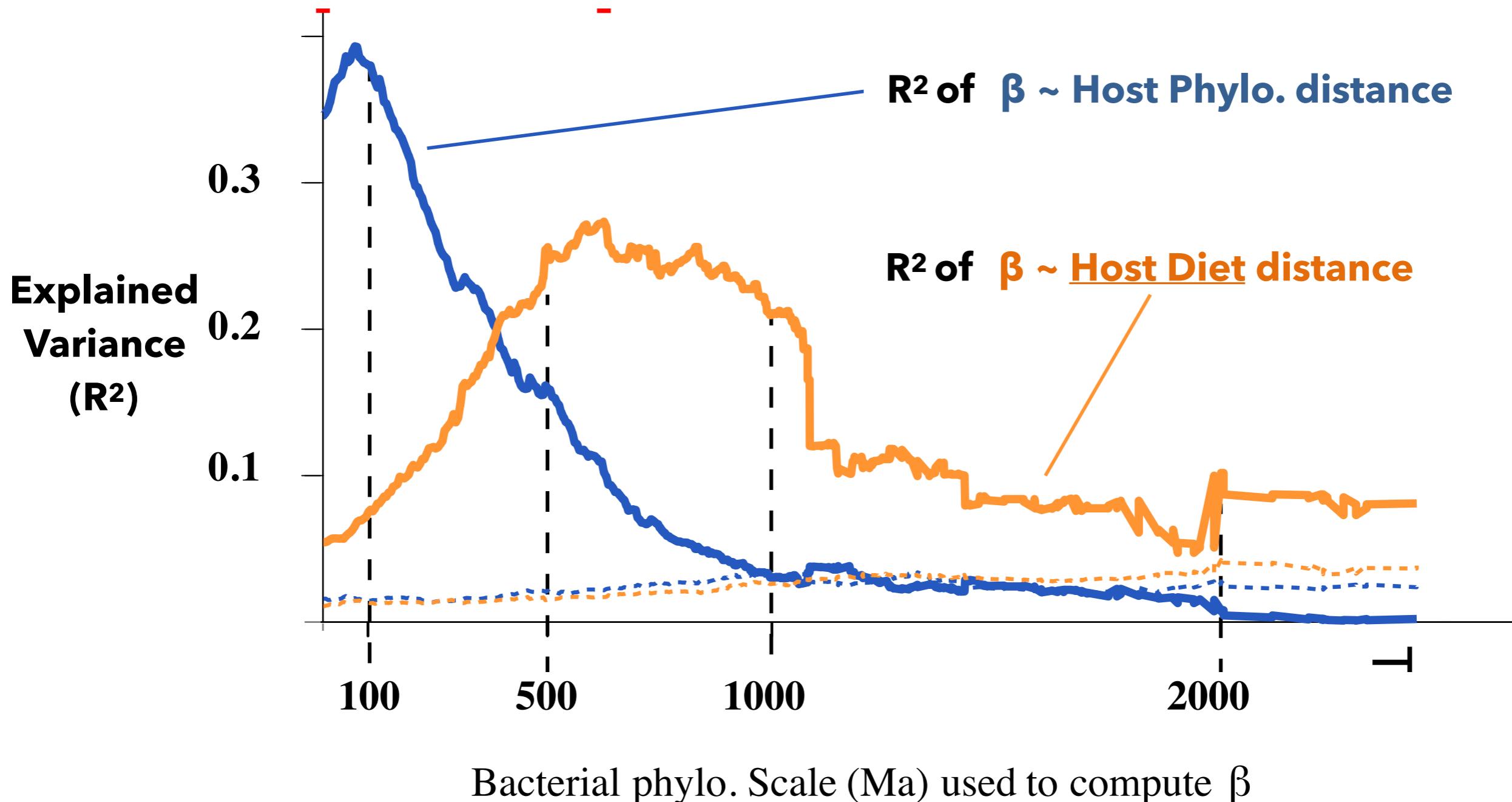


3 – Varying the phylogenetic resolution

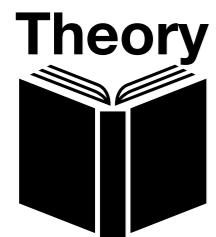


Results: Does that really work ?

Example 1: Mammalian guts

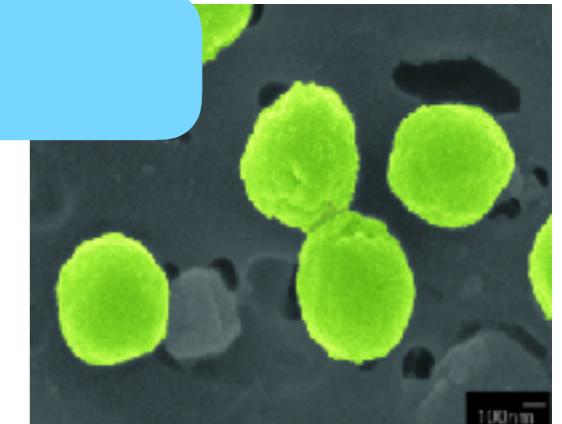


3 – Varying the phylogenetic resolution

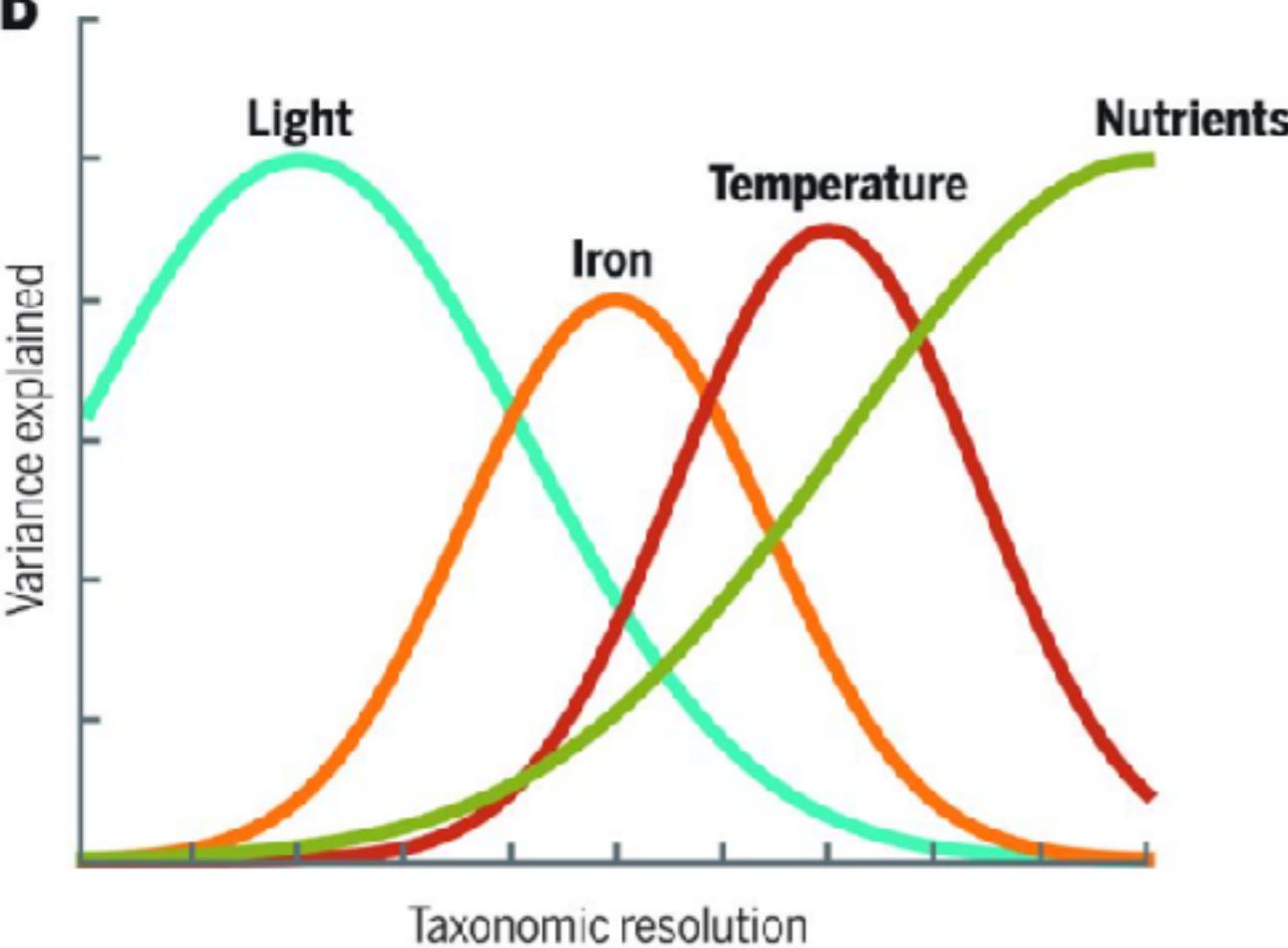


Results: Does that really work ?

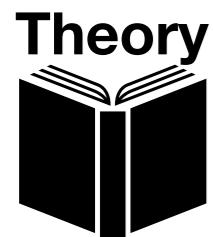
Example 2: *Prochlorococcus* sp.



B



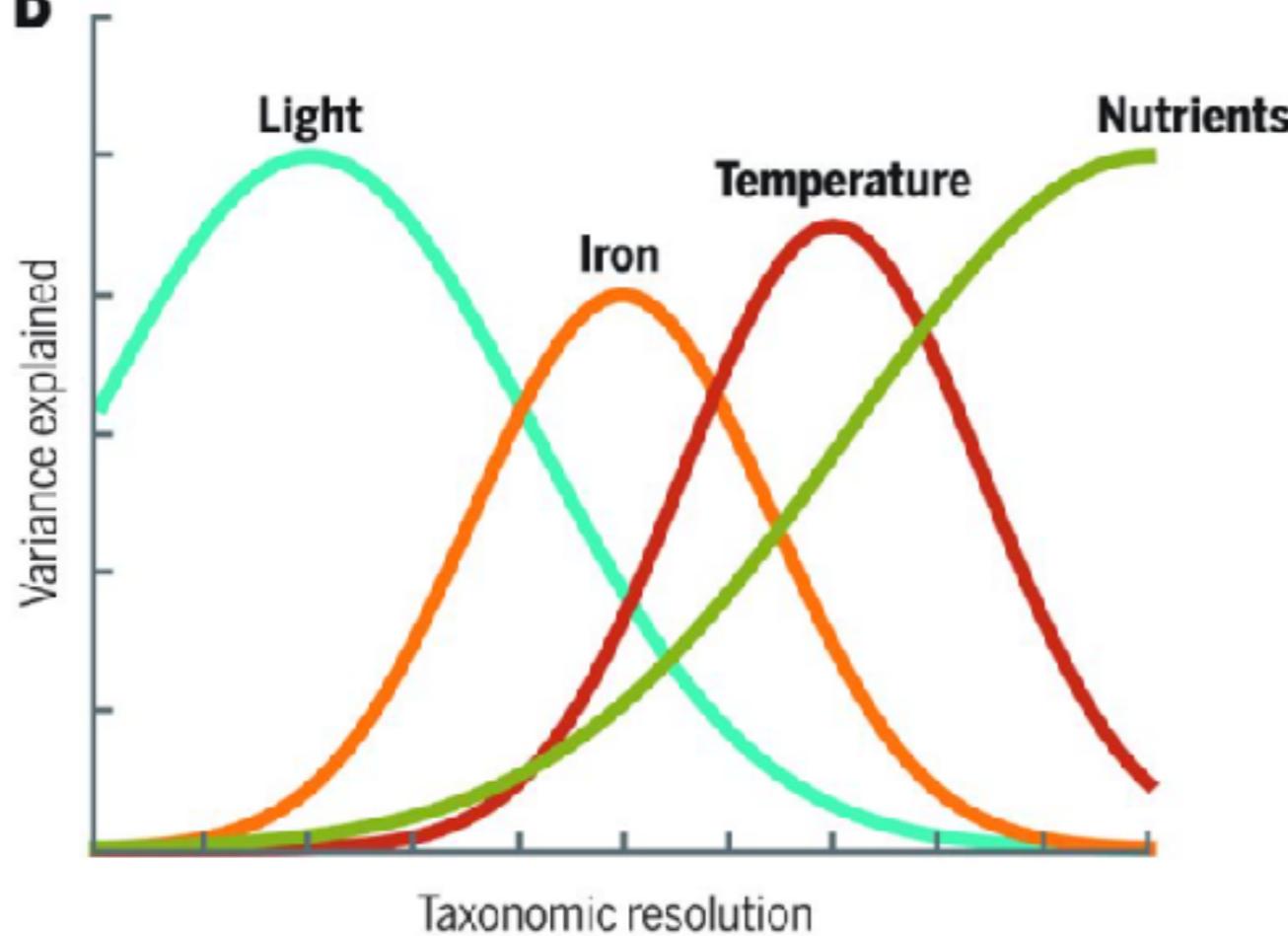
3 – Varying the phylogenetic resolution



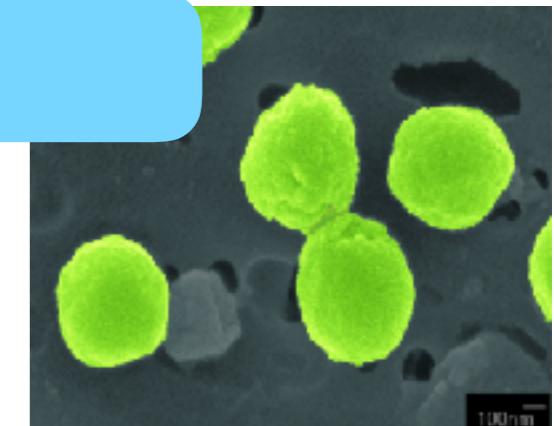
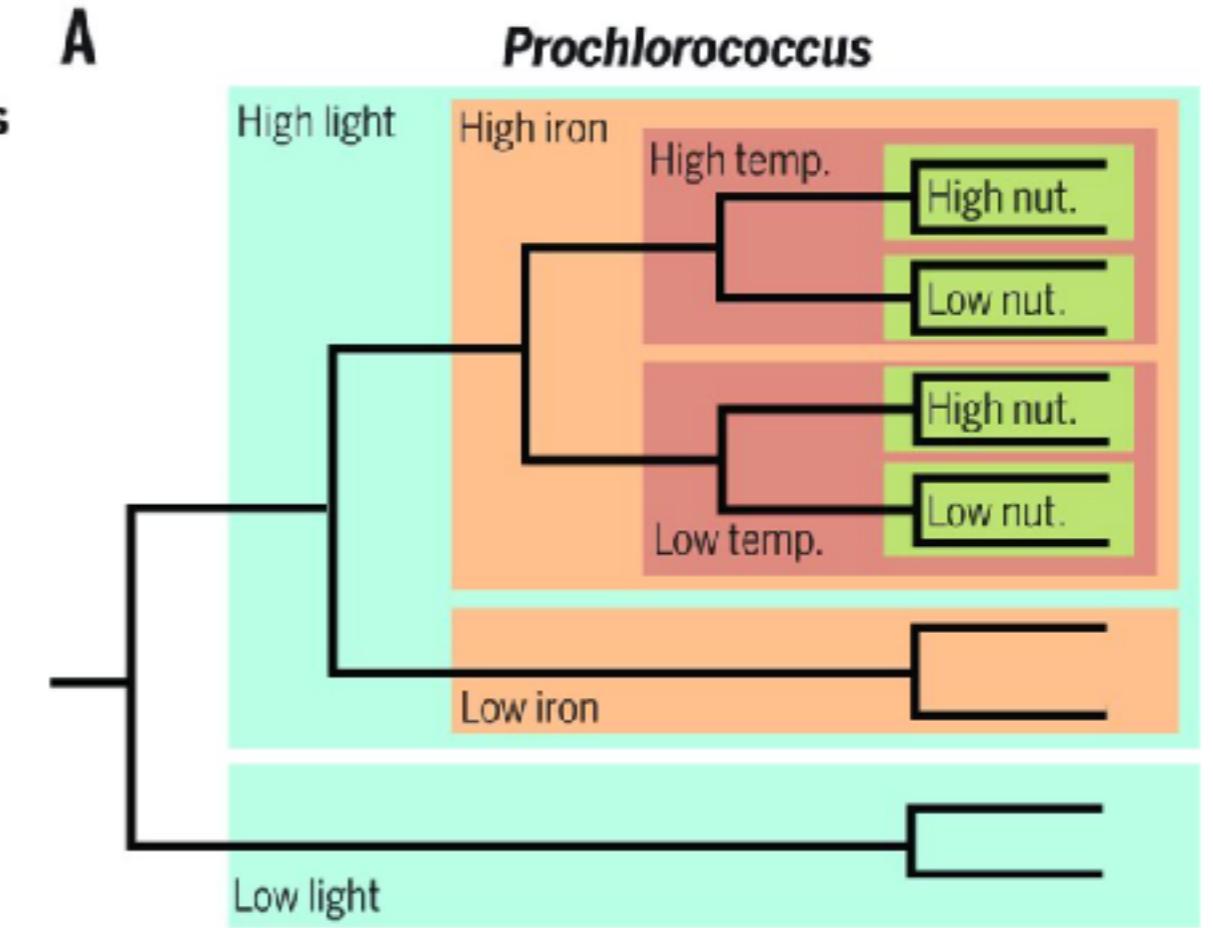
Results: Does that really work ?

Example 2: *Prochlorococcus* sp.

B



A





Get a sense of the phylogenetic scale

`> Hnodes()`

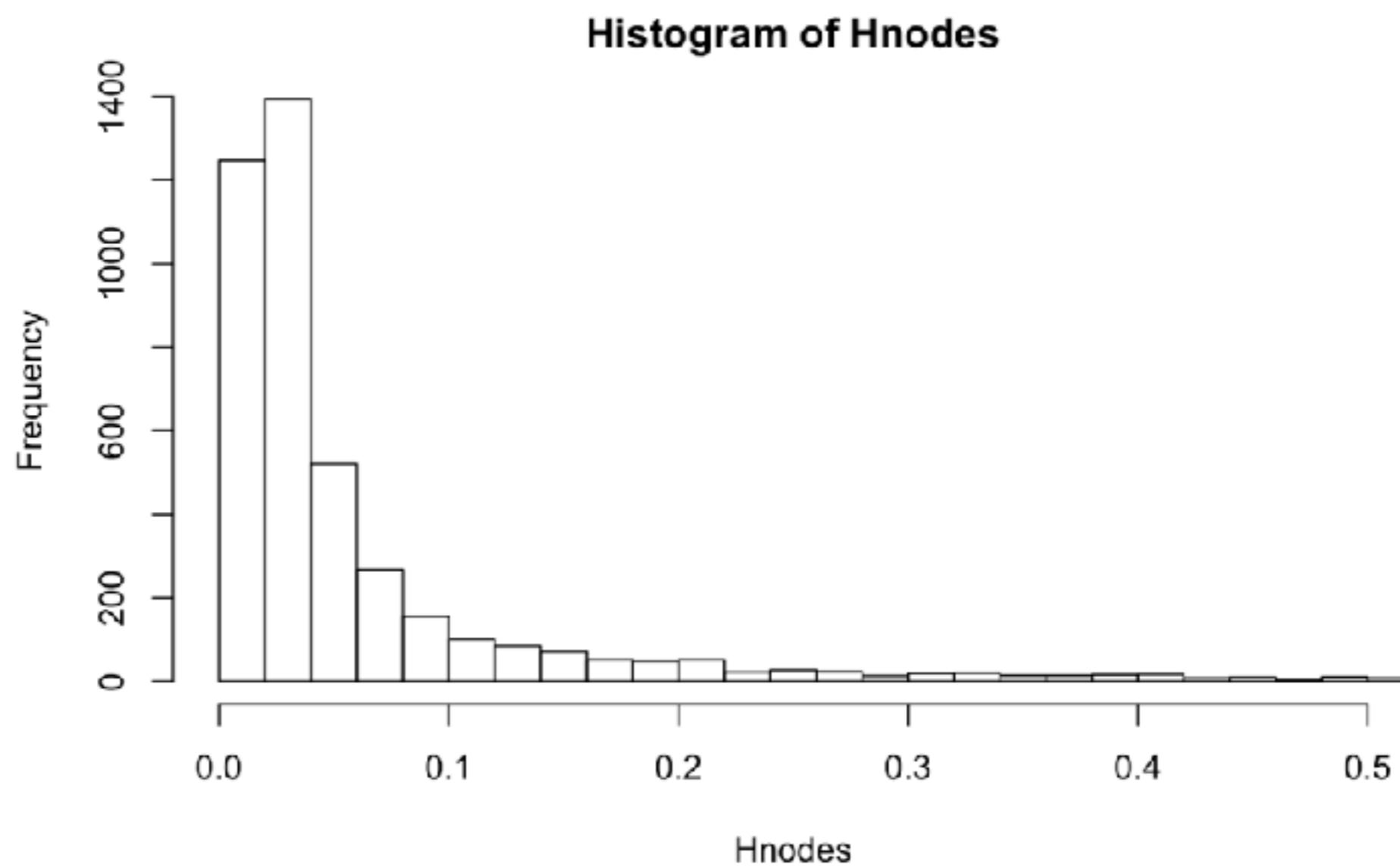
**Compute
Beta diversity through time
(BDTT)**

`> BDTT()`



Get a sense of the phylogenetic scale

> **Hnodes()**





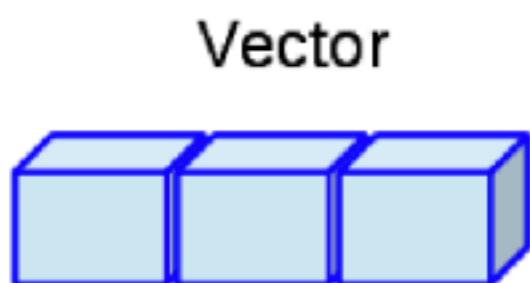
Run the BDTT

> BDTT()

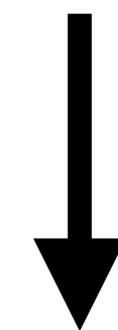
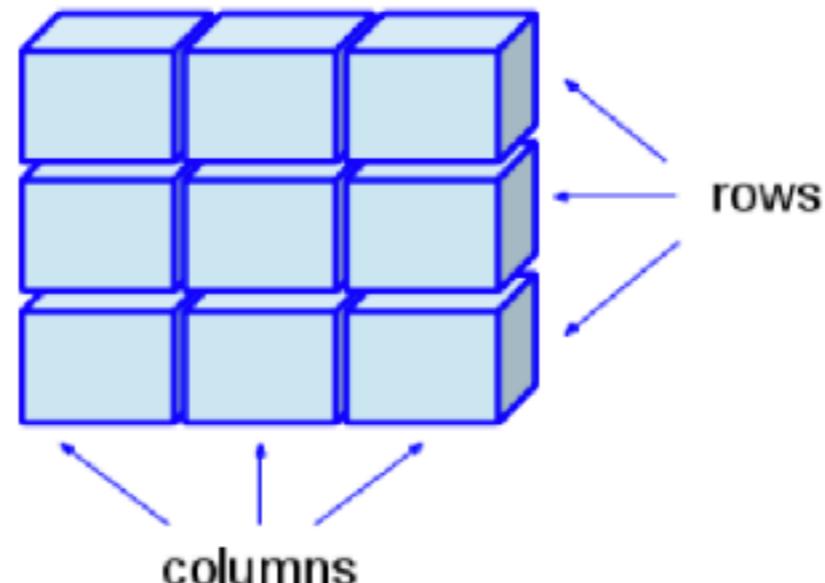
```
> MultipleBetaJac=BDTT(similarity_slices=slices,tree=Tree,sampleOTUs=mat,onlyBeta=T,metric="jac")
[1] "0 similarity provides 4315 total new OTUs"
[1] "0.025 similarity provides 2591 total new OTUs"
[1] "0.05 similarity provides 1365 total new OTUs"
[1] "0.075 similarity provides 932 total new OTUs"
[1] "0.1 similarity provides 733 total new OTUs"
[1] "0.125 similarity provides 606 total new OTUs"
[1] "0.15 similarity provides 515 total new OTUs"
[1] "0.175 similarity provides 435 total new OTUs"
[1] "0.2 similarity provides 378 total new OTUs"
[1] "0.225 similarity provides 322 total new OTUs"
```



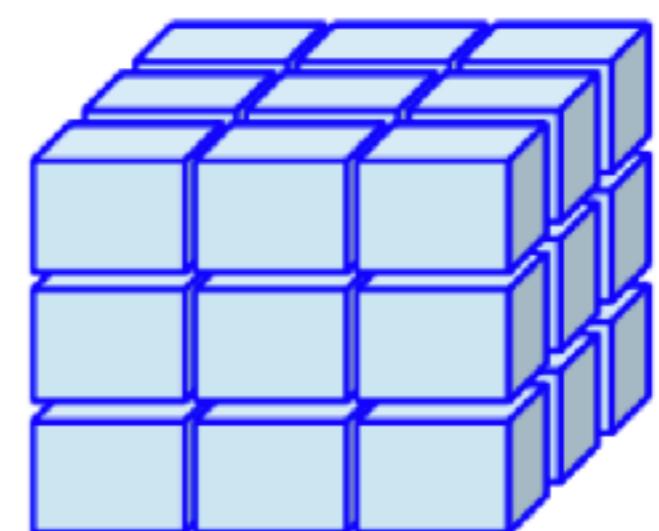
Explore the output



Matrix

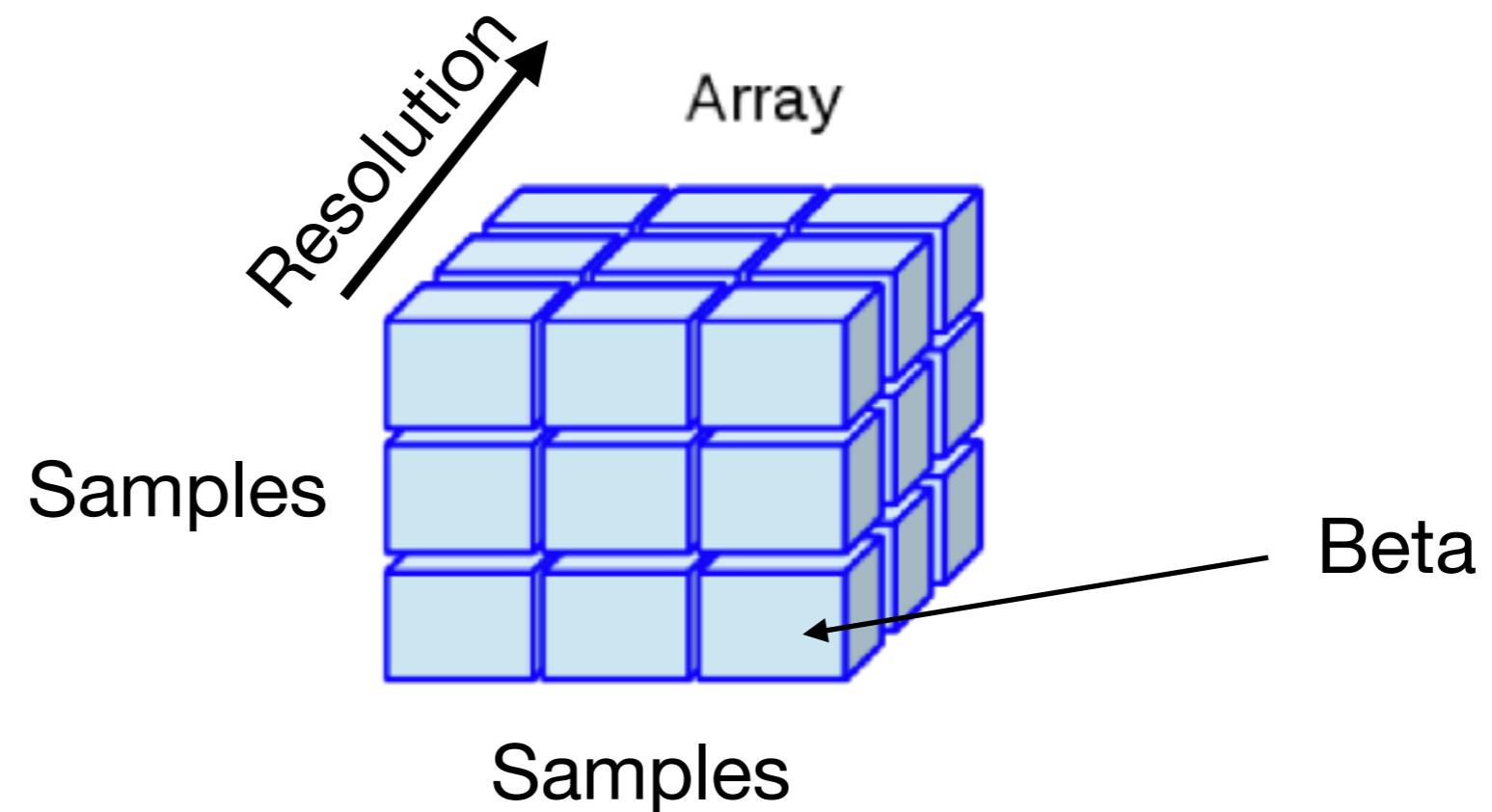


Array



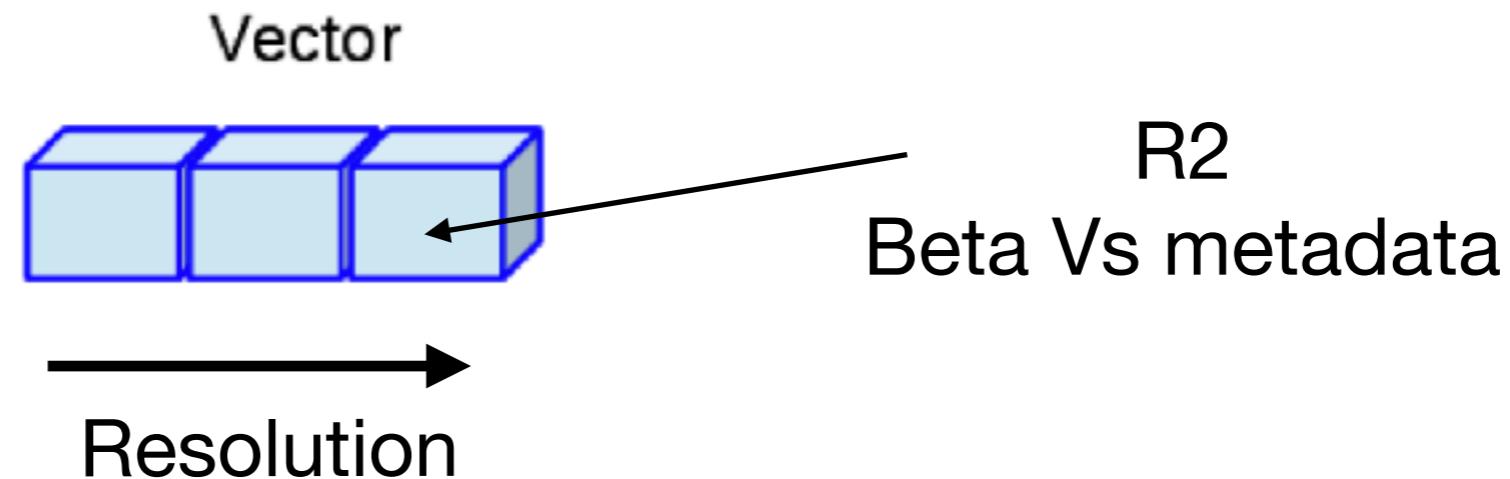


Explore the output

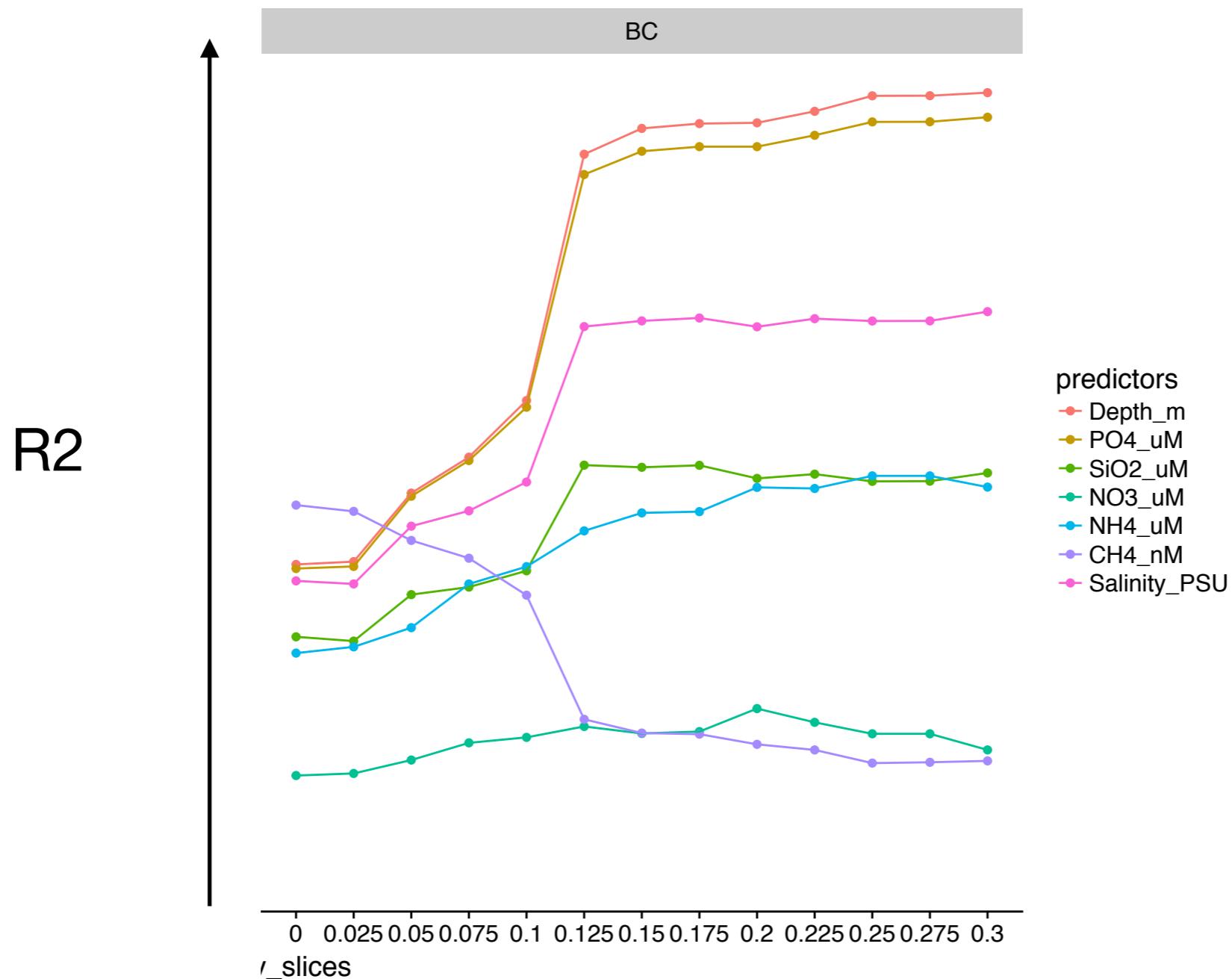




Link with metadata



3 – Varying the phylogenetic resolution



Overall structure of the workshop

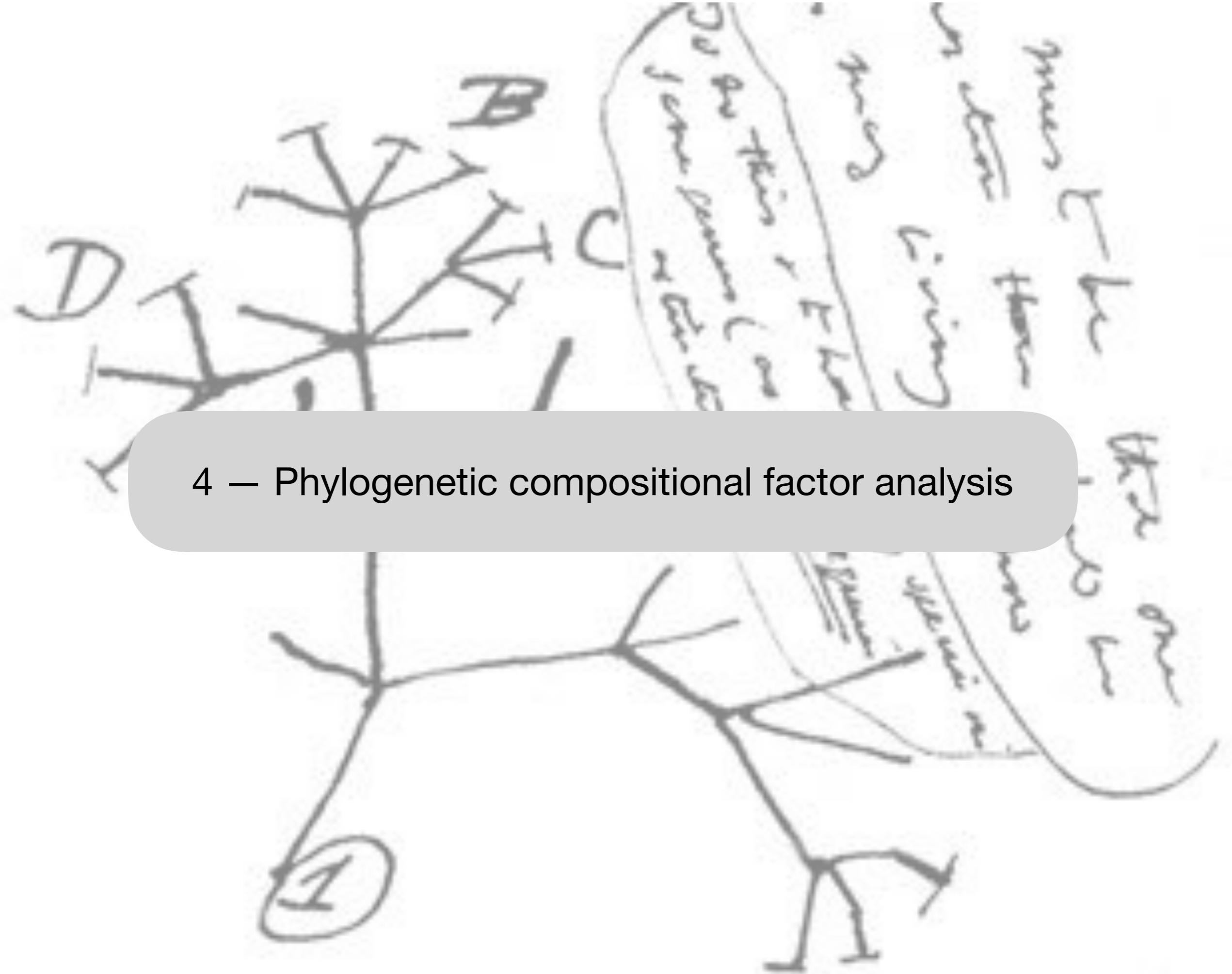
1 – Building a phylogenetic tree

2 – Classical analysis of microbiome

3 – Varying the phylogenetic resolution

4 – Exploring the branches of the phylogenetic tree

4 – Phylogenetic compositional factor analysis



***Two problems:***

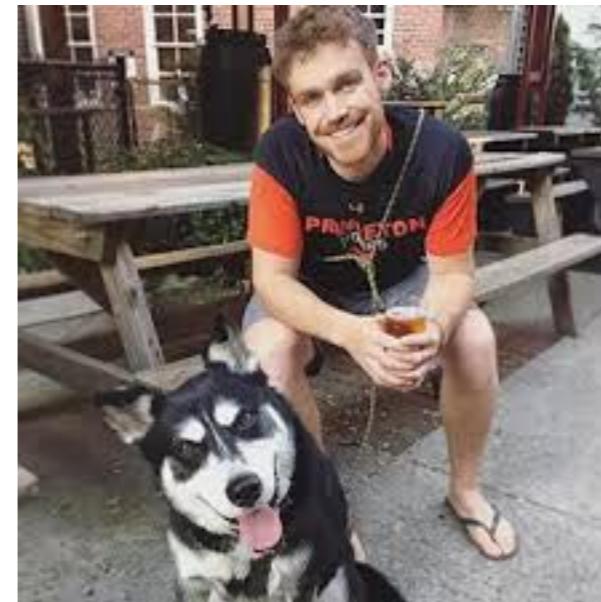
1. *How do you identify changing microbes at unknown levels?*
2. *How do you deal with compositionality?*



4 – PhyloFactor



Jamie Morton

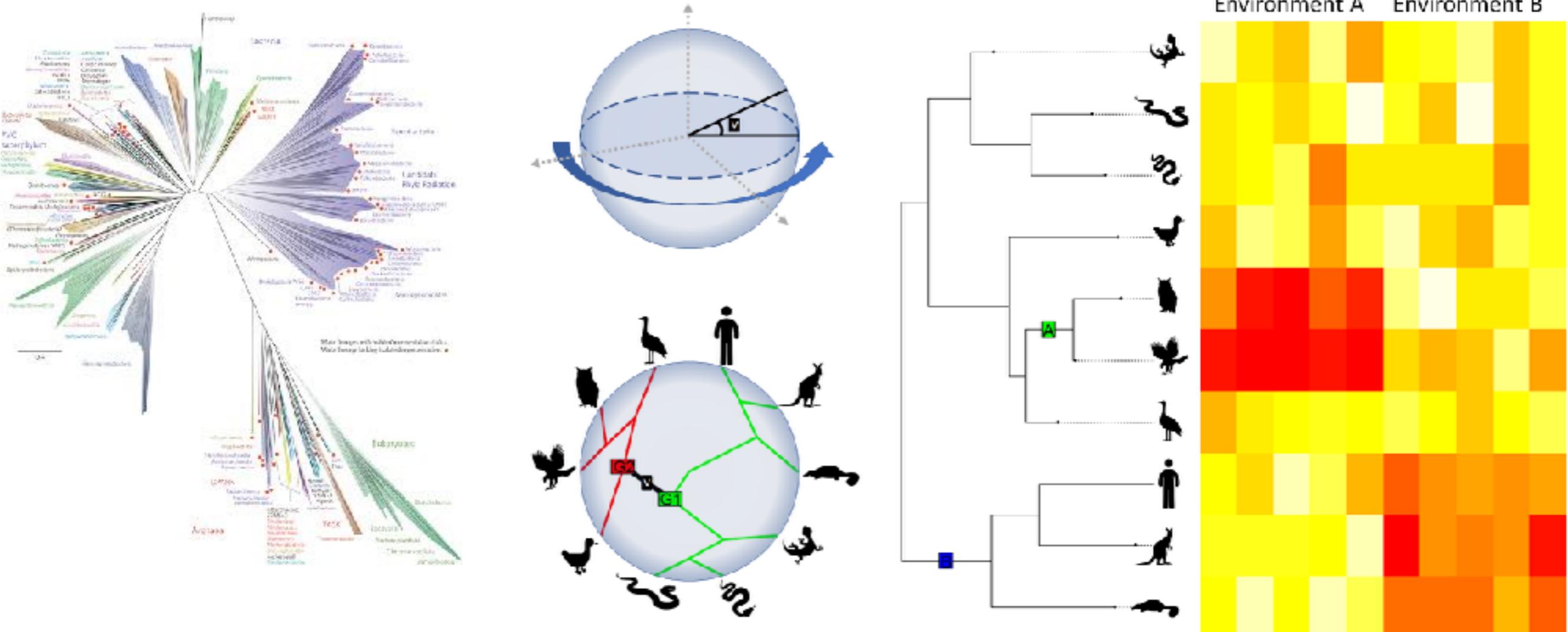


Alex Washburne

Washburne, Morton *et al.* 2018. Methods for phylogenetic analysis of microbiome data. *Nature Microbiology*



Phylogenetic Variables



Washburne, Morton et al. (2018) Methods for phylogenetic analysis of microbiome data. *Nature Microbiology*



Compositional data

For one element to increase, remaining elements must decrease

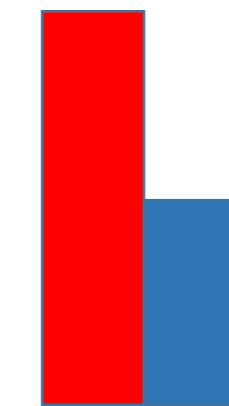


Proportions	1/2	1/2
Counts	100	100



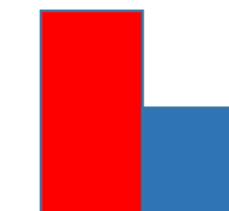
Proportions	2/3	1/3
Counts	200	100

Red doubled



Proportions	2/3	1/3
Counts	100	50

Blue halved



Proportions	1/2	1/2
Counts	100	100

Time point 1

Time point 2



Compositional data

Impossible to know from sequence data alone whether one population is increasing, or another population is decreasing!

*Compositional analyses ask instead: how does
share of community change?*

4 – PhyloFactor

Theory



*Compositional analyses ask instead: how does
share of community change?*



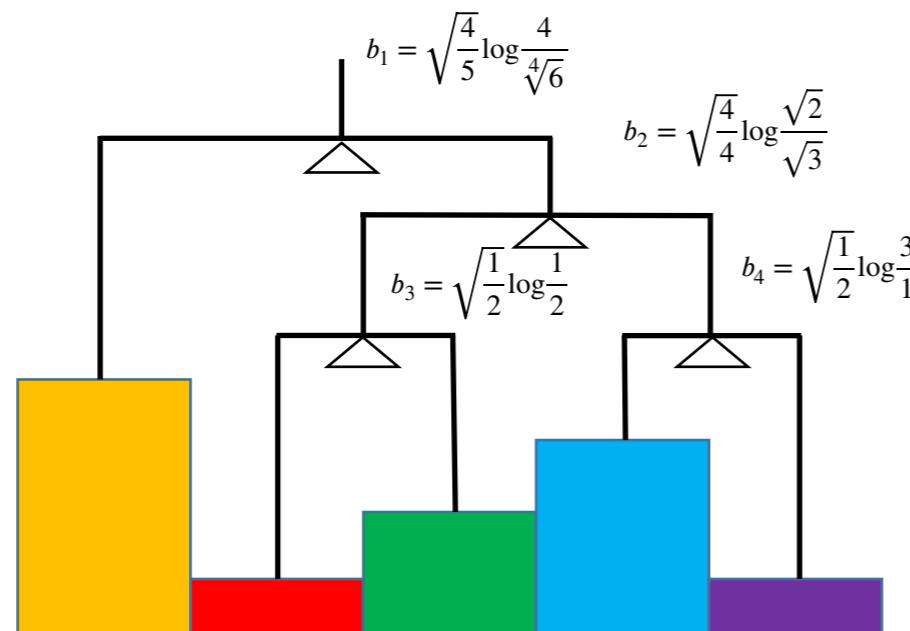


Balance trees

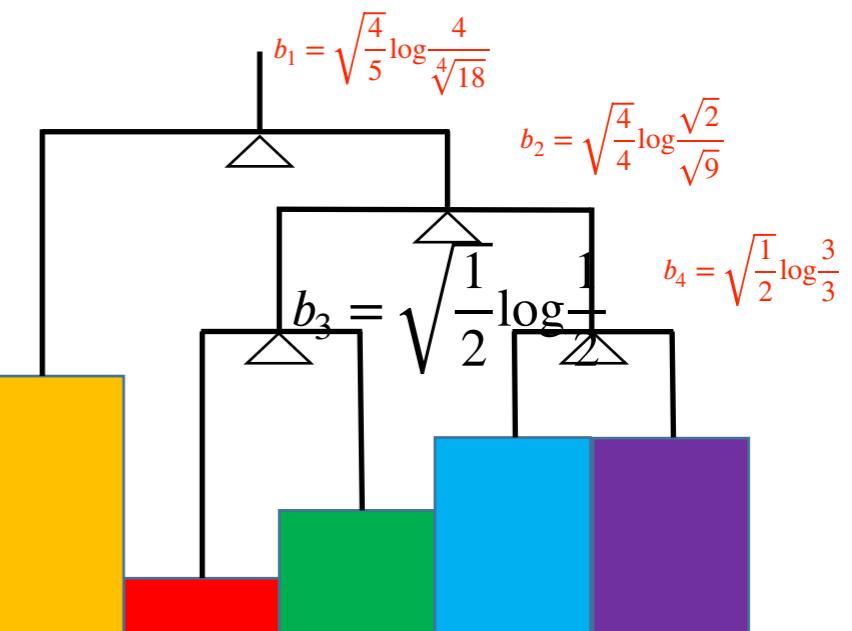
Use a **rooted tree** to construct a *basis* for projecting count data using the **Isometric Log Ratio (ILR)** transform



Balance trees



Proportions	4/11	1/11	2/11	3/11	1/11
Percentages	36%	9%	18%	27%	9%
Counts	4	1	2	3	1



Proportions	4/13	1/13	2/13	3/13	3/13
Percentages	30%	7%	14%	21%	21%
Counts	4	1	2	3	3

$$b_i = \sqrt{\frac{|i_L| |i_R|}{|i_L| + |i_R|}} \log \left(\frac{g(i_L)}{g(i_R)} \right)$$

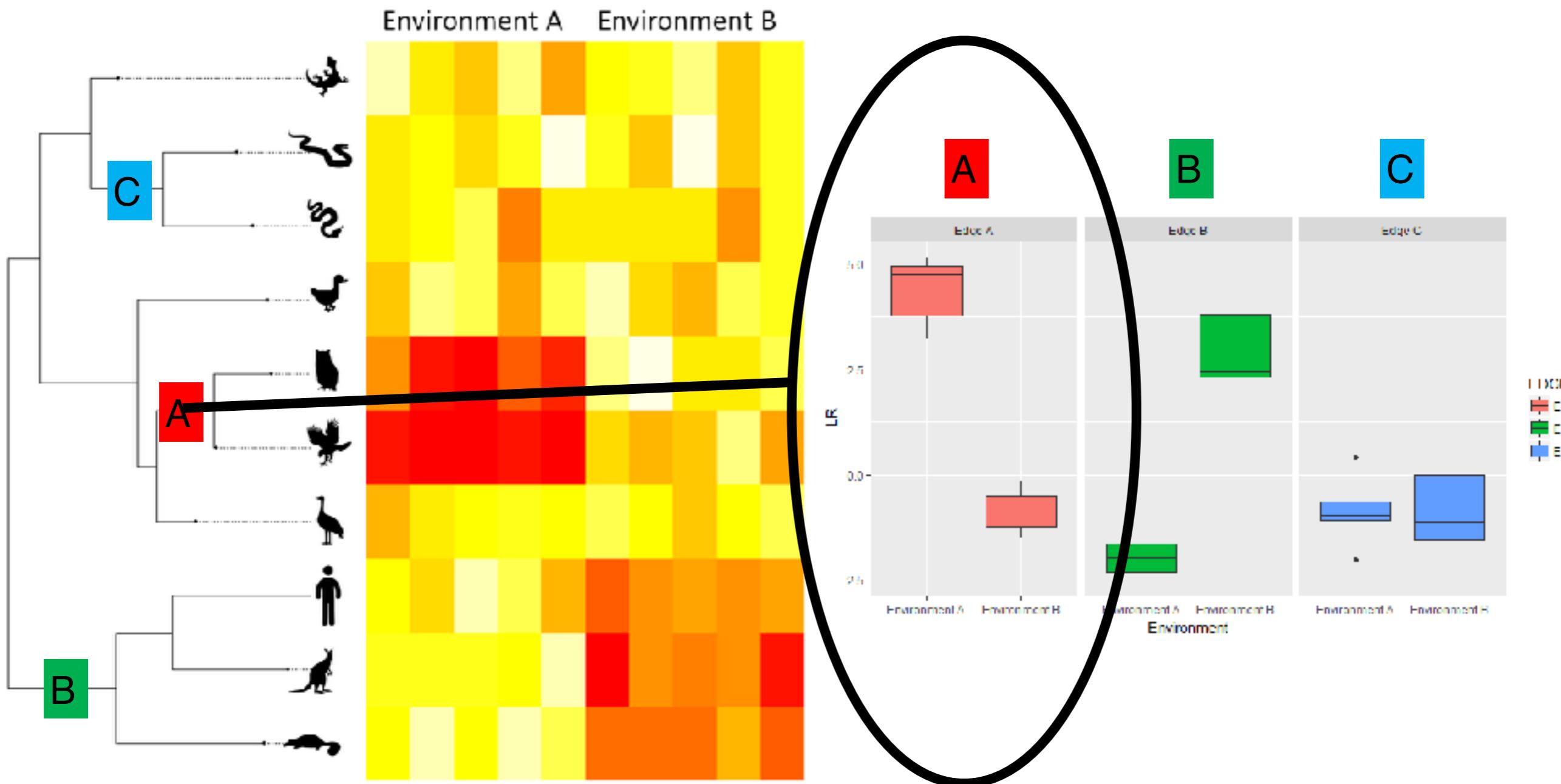


PhyloFactorization

Uses an **unrooted tree** to iteratively construct a ILR basis, identifying ‘most important’ edge on tree (balance) at each iteration

4 — PhyloFactor

Theory



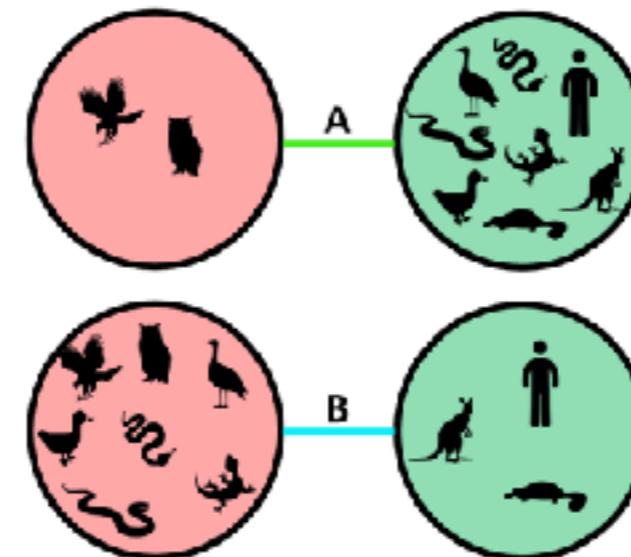
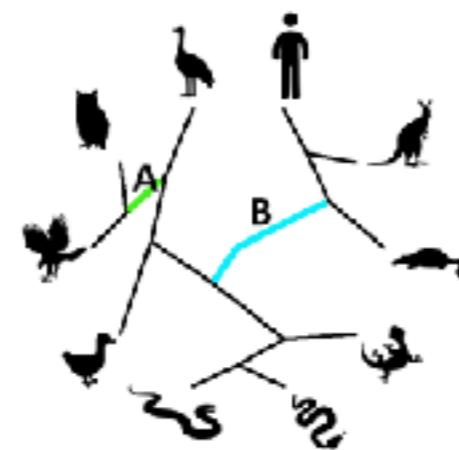
Washburne et al.
(2017) *PeerJ*

4 – PhyloFactor

Theory



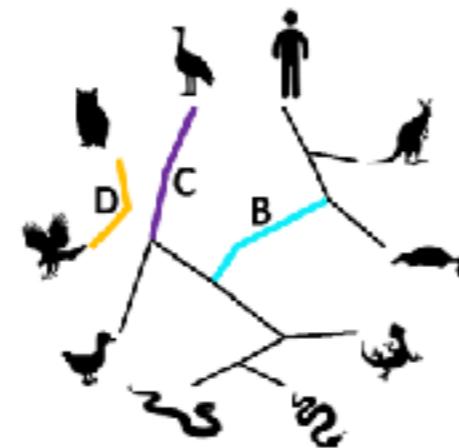
1st Iteration



$$\omega(A) = 2$$

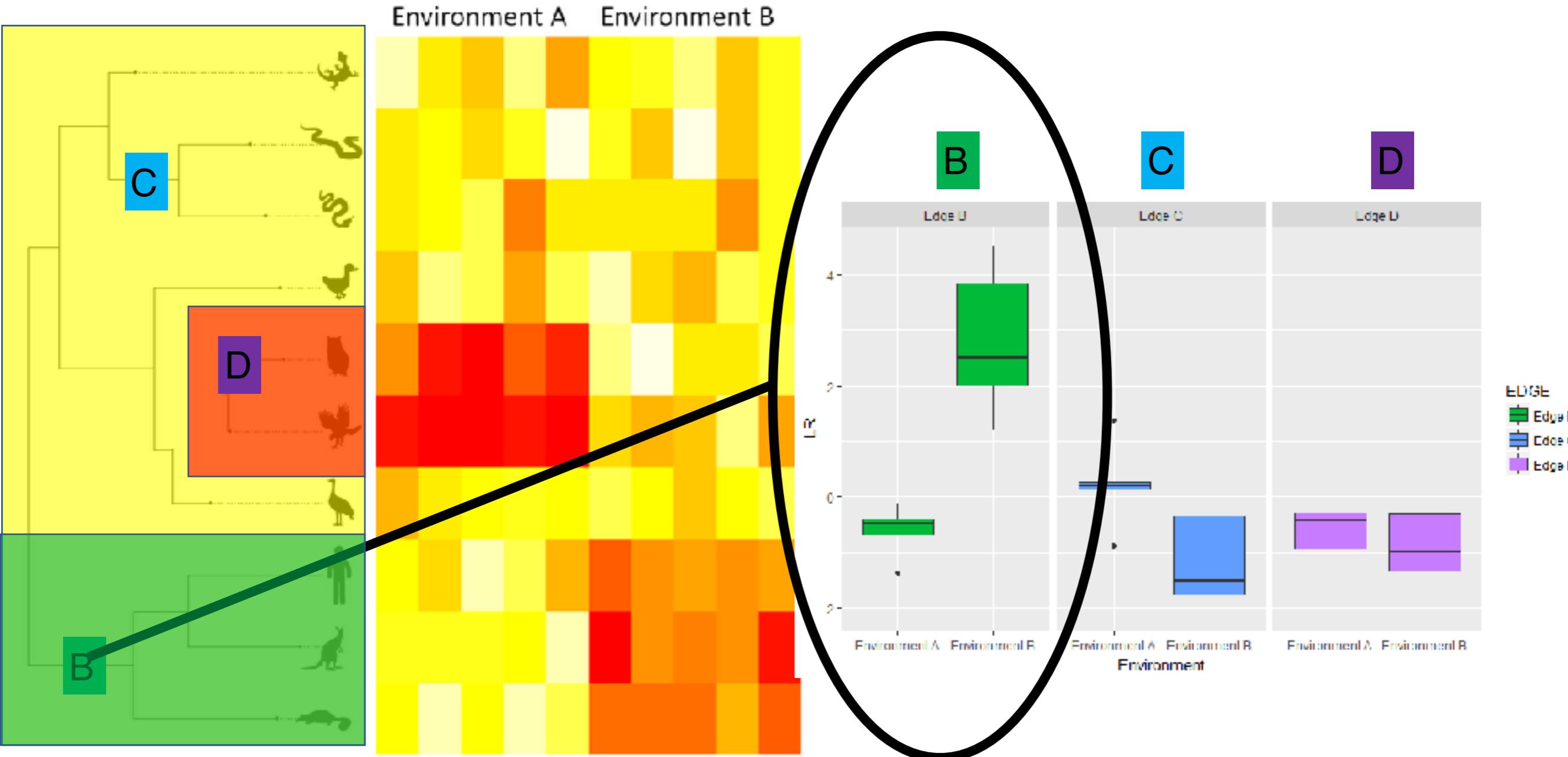
$$\omega(B) = 1$$

2nd Iteration



4 — PhyloFactor

Theory



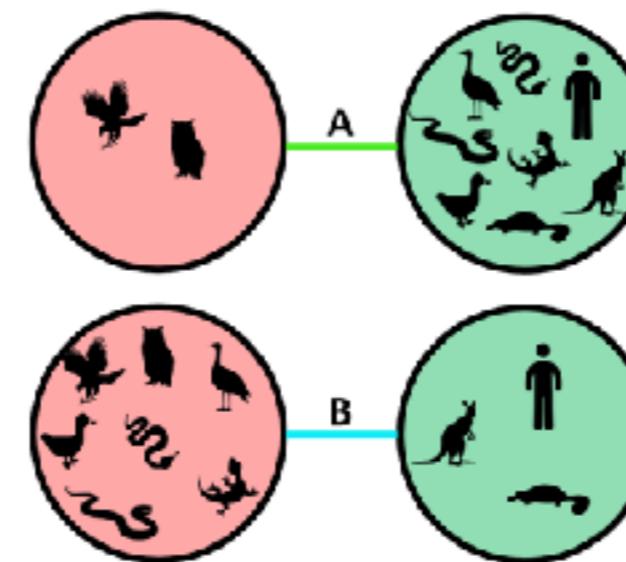
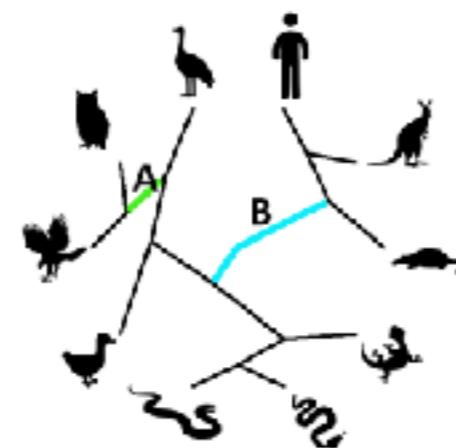
Washburne et al.
(2017) *PeerJ*

4 – PhyloFactor

Theory



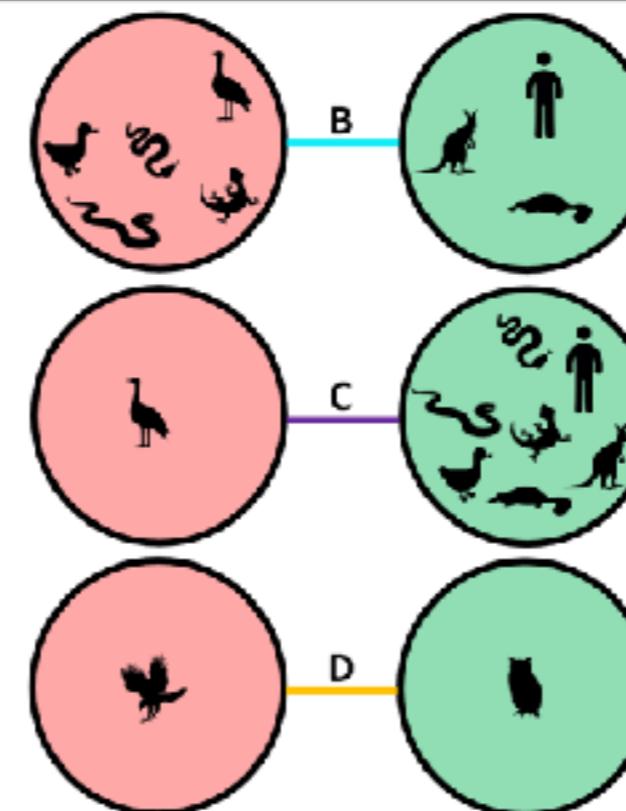
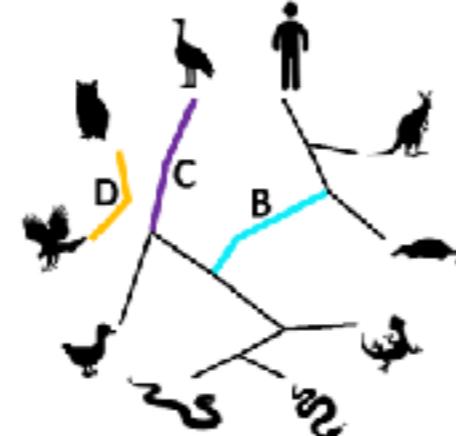
1st Iteration



$$\omega(A) = 2$$

$$\omega(B) = 1$$

2nd Iteration



$$\omega(B) = 3$$

$$\omega(C) = 1$$

$$\omega(D) = \frac{1}{2}$$

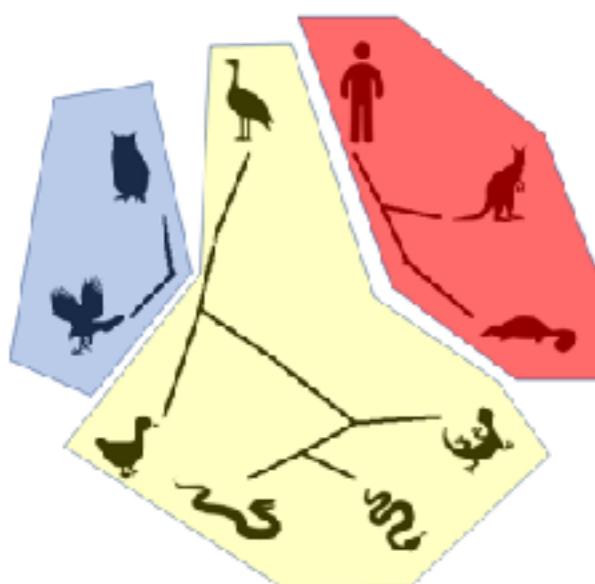
4 — PhyloFactor

Theory

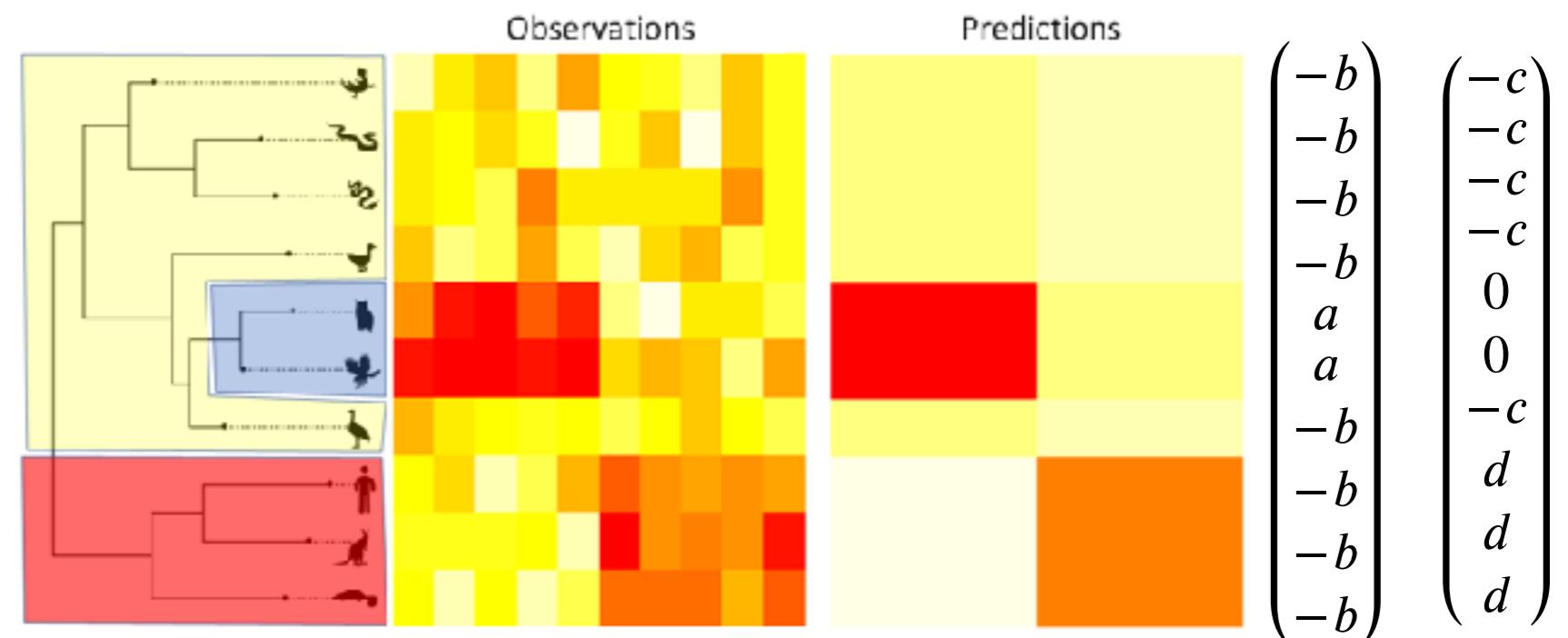


Results

Bins



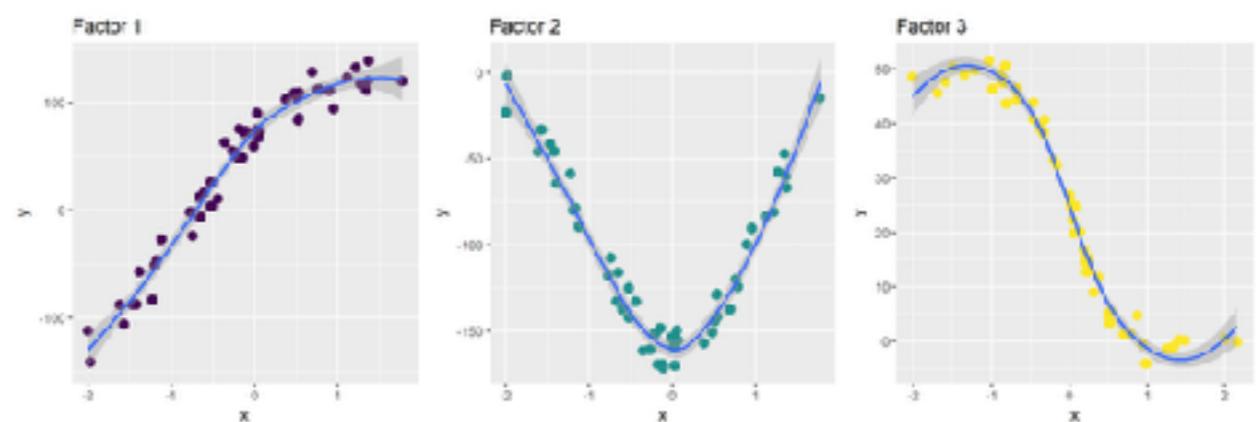
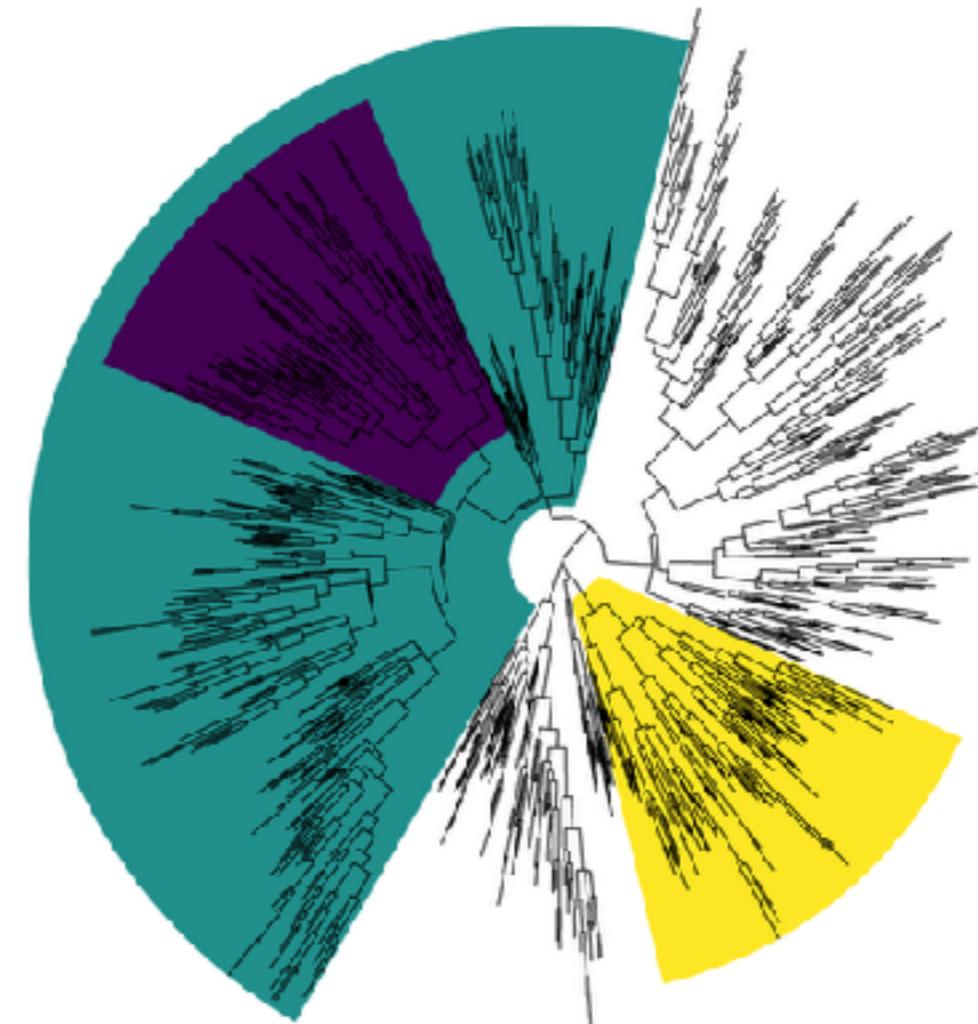
Regression-Phylofactorization





PhyloFactorization

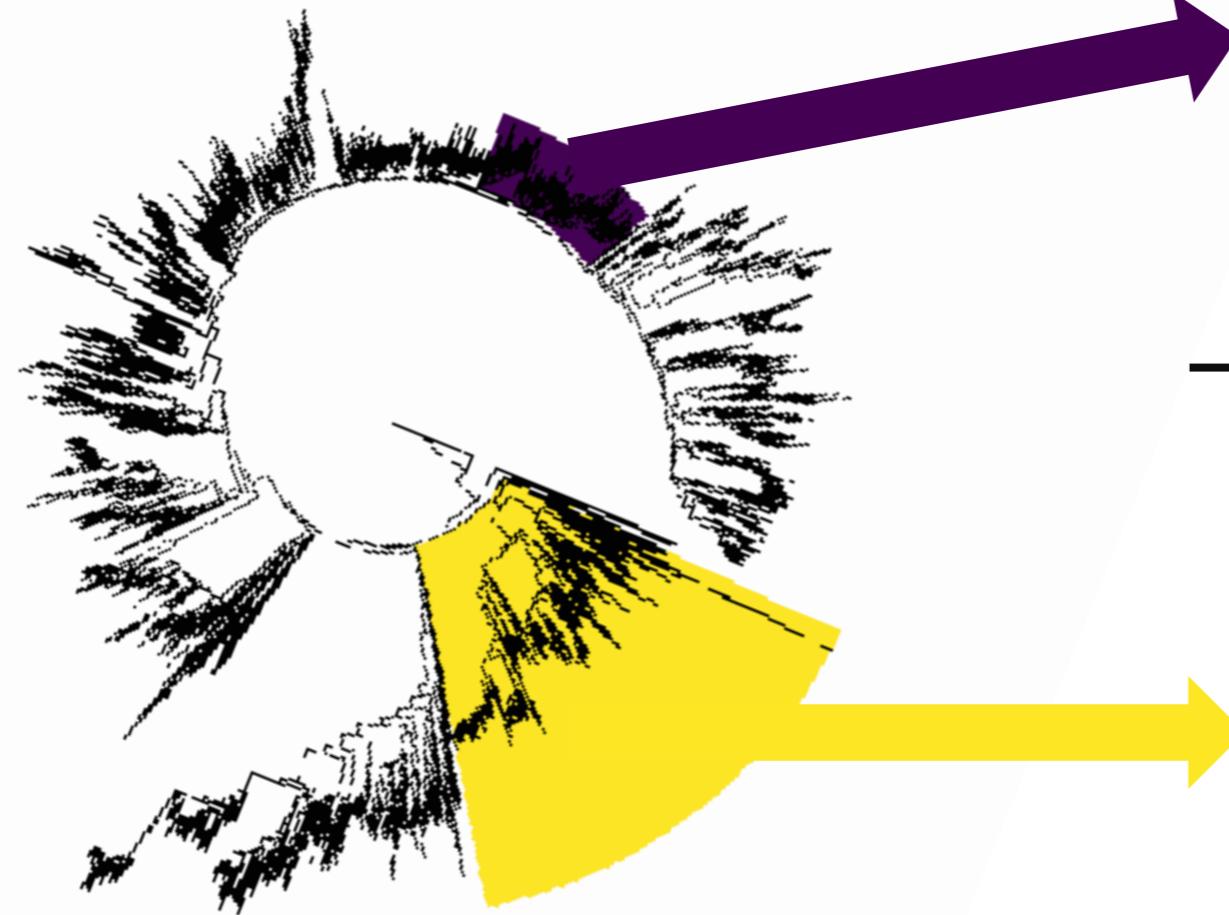
ILR balances are **independent** and **~normal**: allow parametric modeling of phylogenetic variables with more complex models



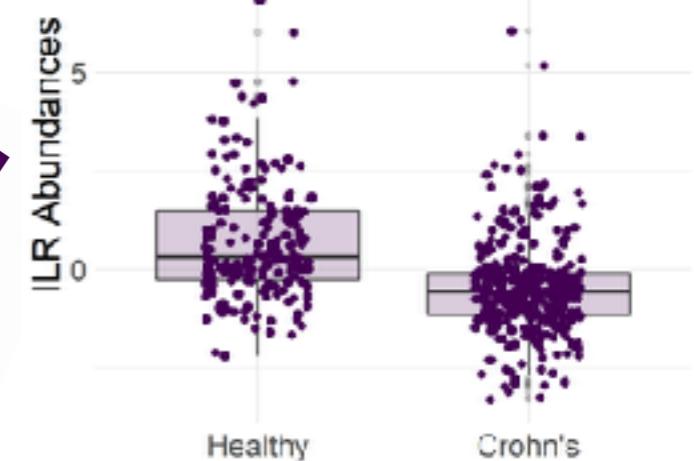


Example 1: Chron's Disease

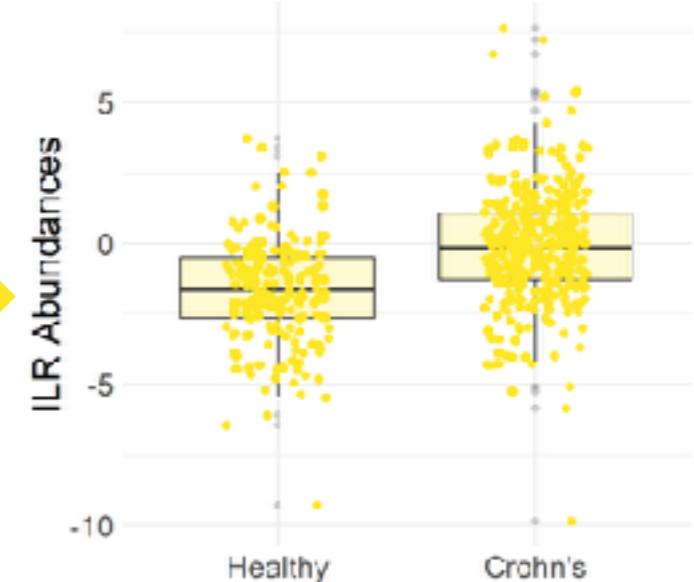
Formula: CD~Data



Lachnospiraceae genera



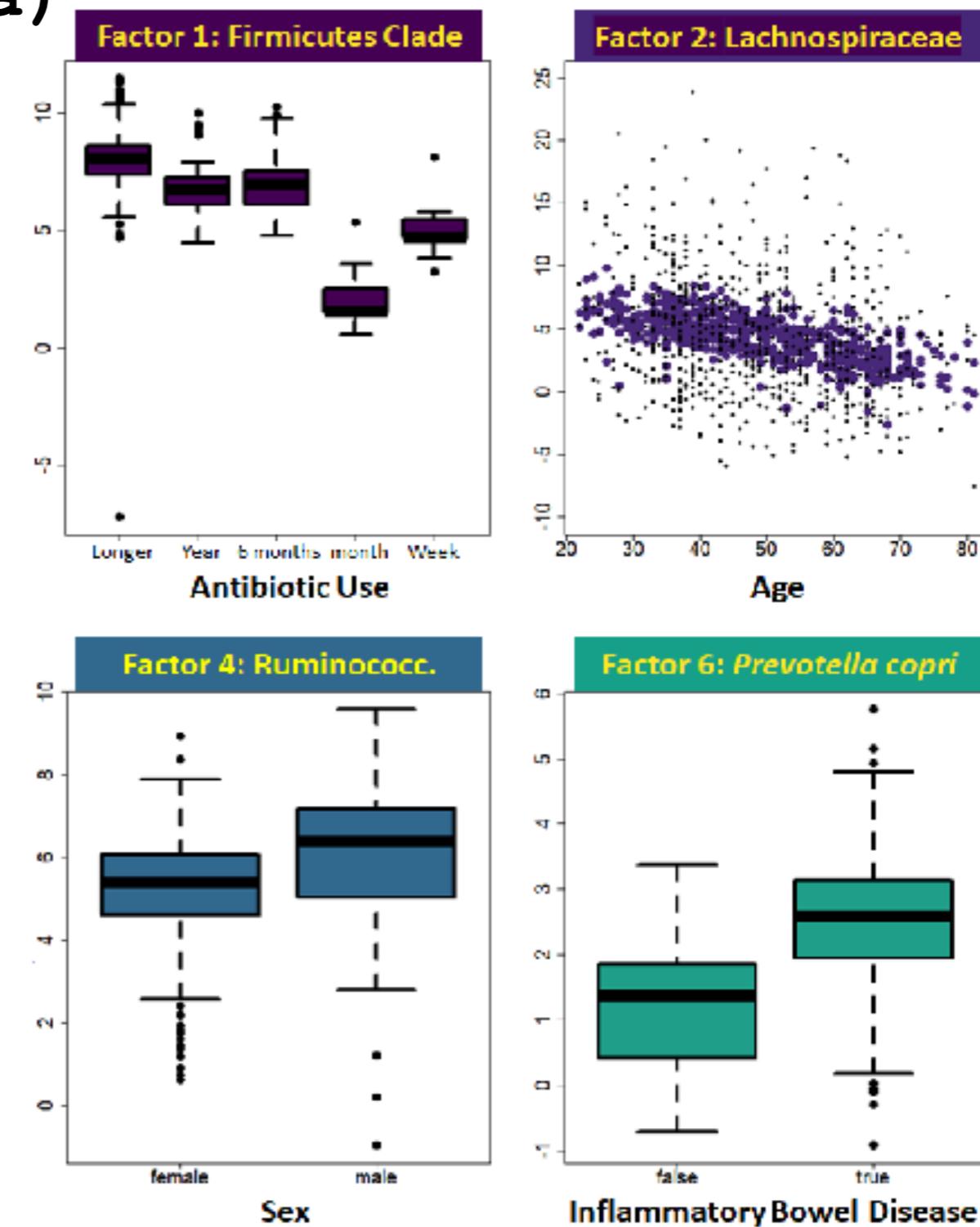
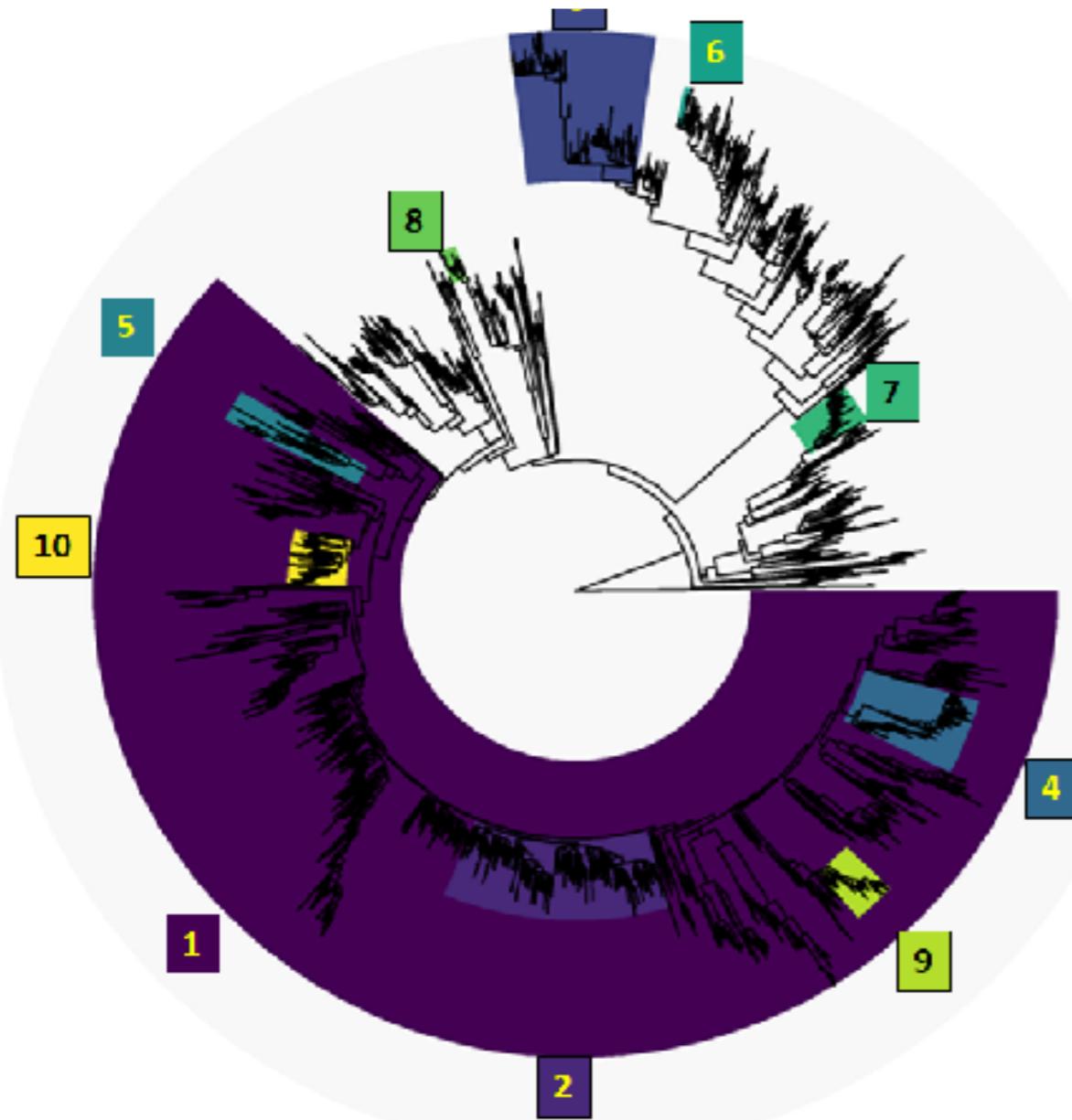
Proteobacterial classes





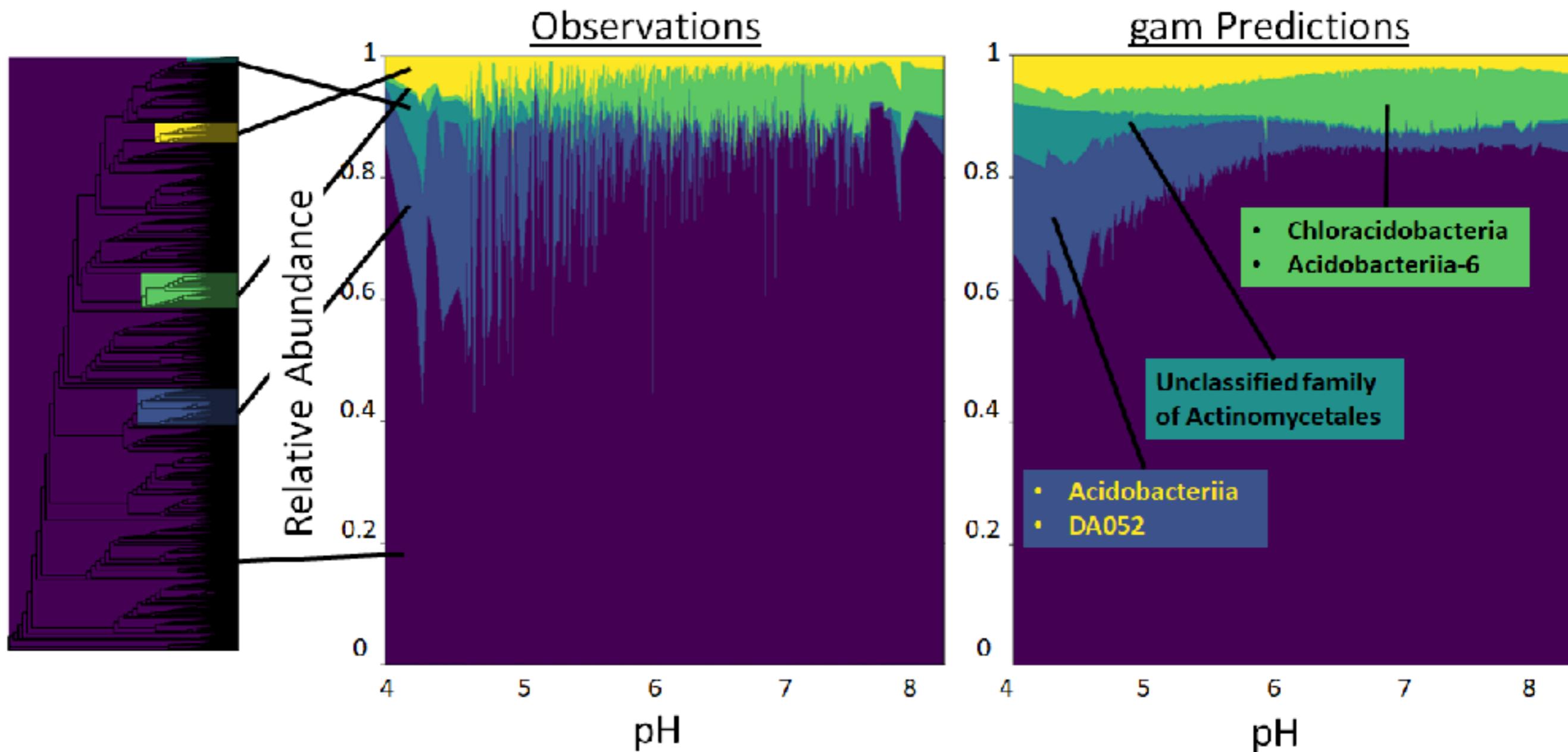
Example 2: American Gut

Formula: `var (Data)`



**Example 3: Central Park Soils (gam)**

Formula: $\text{Data} \sim s(\text{C}) + s(\text{pH}) + s(\text{N})$





Code lines 427 – 724

Import data

```
> read.table()  
> read.tree()
```

Simulating phylo variable

```
> rlnorm()
```

Testing phylo variable with
2-sample test

```
> twoSampleFactor()
```

Testing distribution of OTUs

```
> phylofactor()
```

Summarizing factor taxa

```
> pf.taxa()
```

Visualizing factors on trees

```
> pf.tree()  
> pf.heatmap()
```

Predicting data on trees

```
> predict()
```

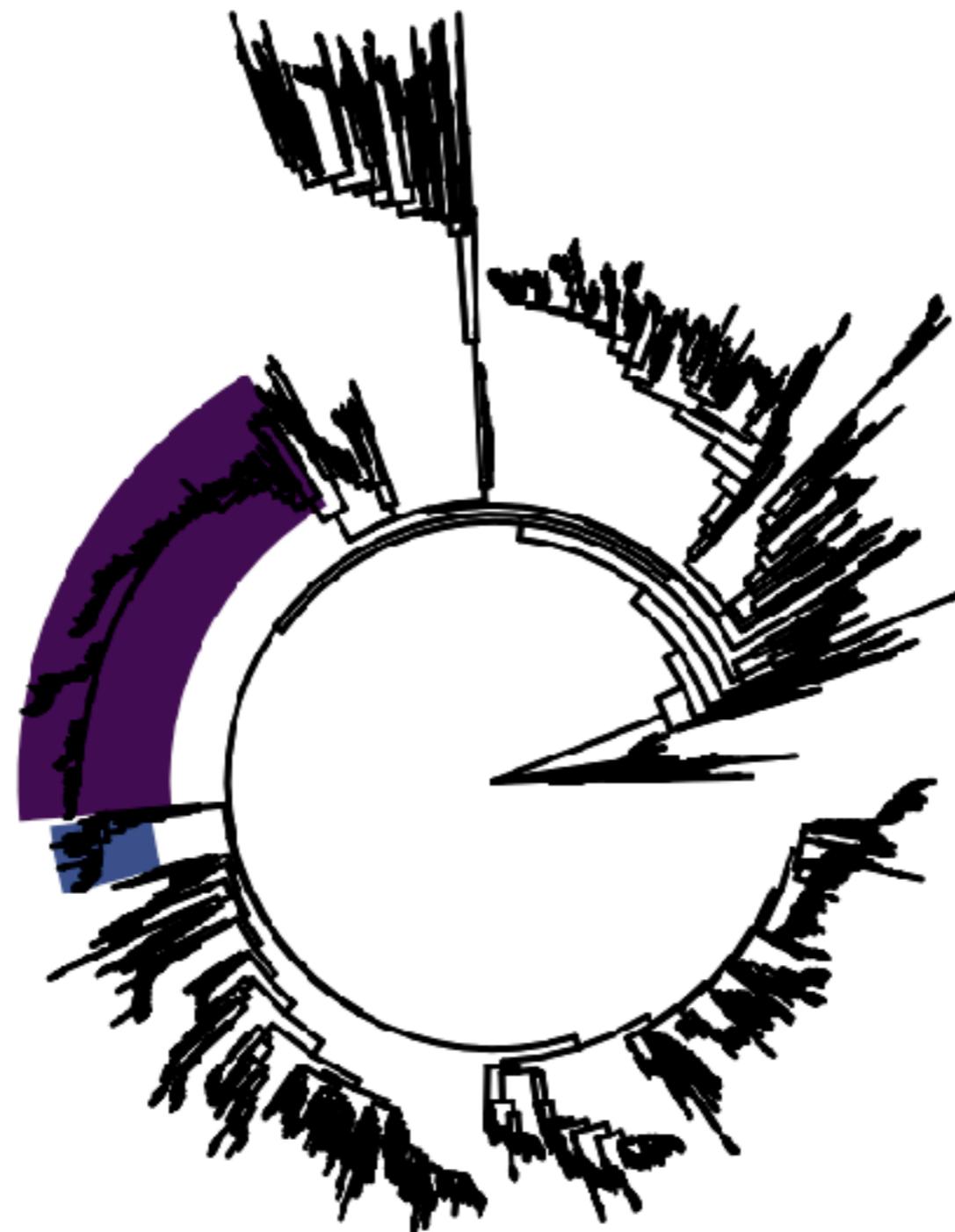


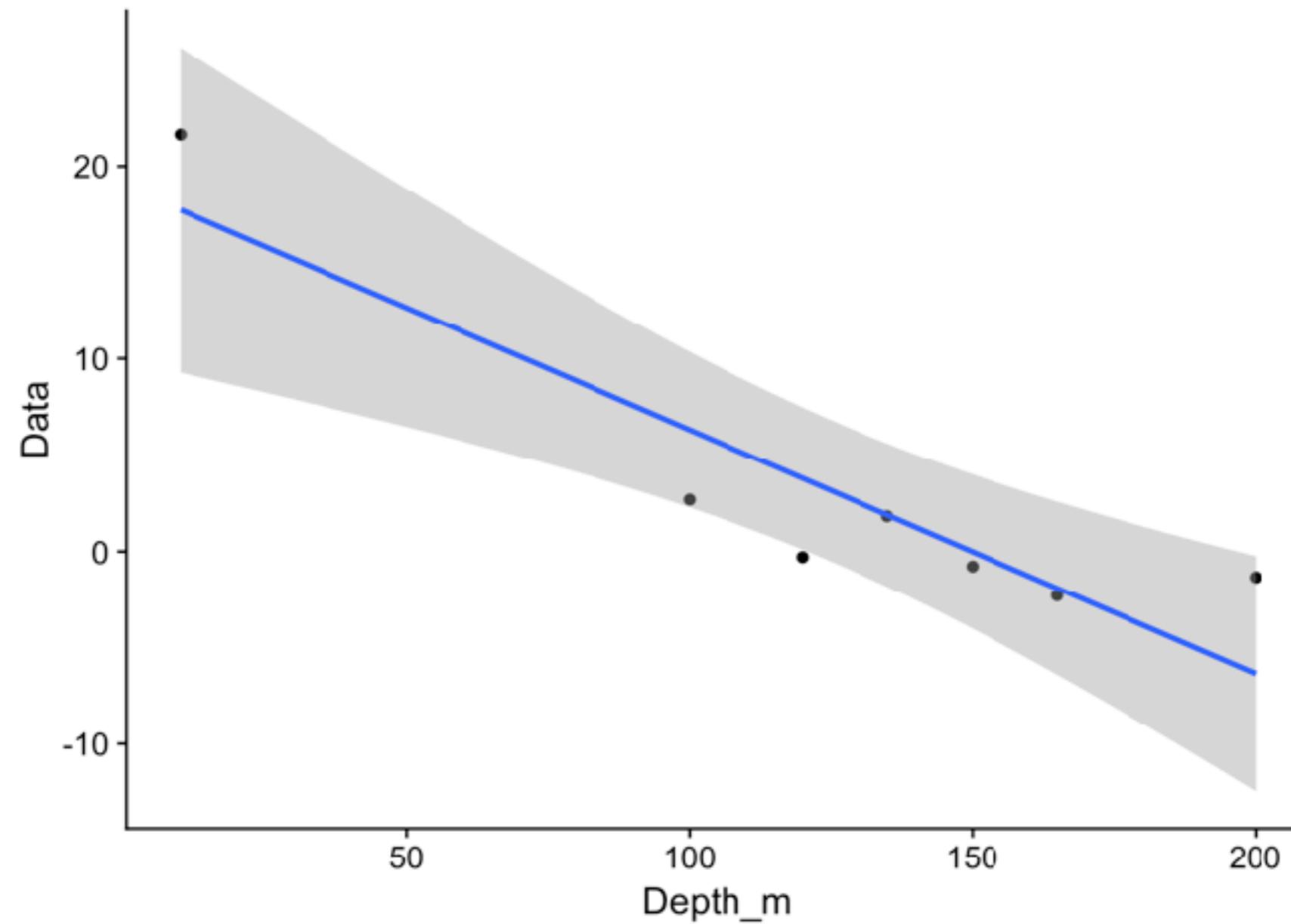
```

##      phylofactor object from function PhyloFactor
## -----
## Method          : glm
## Choice          : var
## Formula         : Data ~ Depth_m
## Number of species : 4315
## Number of factors  : 5
## Frac Explained Variance : 0.0624
## Largest non-remainder bin : 473
## Number of singletons   : 0
## Paraphyletic Remainder : 3250 species
##
## -----
## Factor Table:
##                               Group1          Group2
## Factor 1 473 member Monophyletic clade 3842 member Monophyletic clade
## Factor 2 141 member Monophyletic clade 3701 member Paraphyletic clade
## Factor 3 337 member Monophyletic clade 3364 member Paraphyletic clade
## Factor 4 114 member Monophyletic clade 3250 member Paraphyletic clade
## Factor 5  38 member Monophyletic clade  435 member Paraphyletic clade
##             ExpVar      F     Pr(>F)
## Factor 1 0.0178330 12.333 0.0170660
## Factor 2 0.0129200 24.596 0.0042502
## Factor 3 0.0136960 45.169 0.0011049
## Factor 4 0.0105830 42.647 0.0012596
## Factor 5 0.0073844 22.799 0.0049935

## [1] "Bacteria;Proteobacteria;Alphaproteobacteria;Rhizobiales;PS1_clade"
## [2]
## "Bacteria;Proteobacteria;Alphaproteobacteria;Alphaproteobacteria_unclassified;Alphaproteobacteria_unclassified"
## [3] "Bacteria;Proteobacteria;Alphaproteobacteria;Rhizobiales;Rhizobiales_unclassified;Rhizobiales_unclassified"
## [4] "Bacteria;Proteobacteria;Alphaproteobacteria;Rhodobacterales;Rhodobacteraceae;Rhodobacteraceae"
## [5] "Bacteria;Proteobacteria;Alphaproteobacteria;Rhodobacterales;Rhodobacteraceae;Frigidibacter"
## [6] "Bacteria;Proteobacteria;Alphaproteobacteria;Rhizobiales;Phyllobacteriaceae"
## [7] "Bacteria;Proteobacteria;Alphaproteobacteria;Caulobacterales;Hyphomonadaceae;Hellea"
## [8] "Bacteria;Proteobacteria;Alphaproteobacteria;Rhodobacterales;Rhodobacteraceae;Octadecabacter"
## [9] "Bacteria;Proteobacteria;Alphaproteobacteria;Rhodobacterales;Rhodobacteraceae;Litoreibacter"
## [10] "Bacteria;Proteobacteria;Alphaproteobacteria;Rhodobacterales;Rhodobacteraceae;Amylibacter"

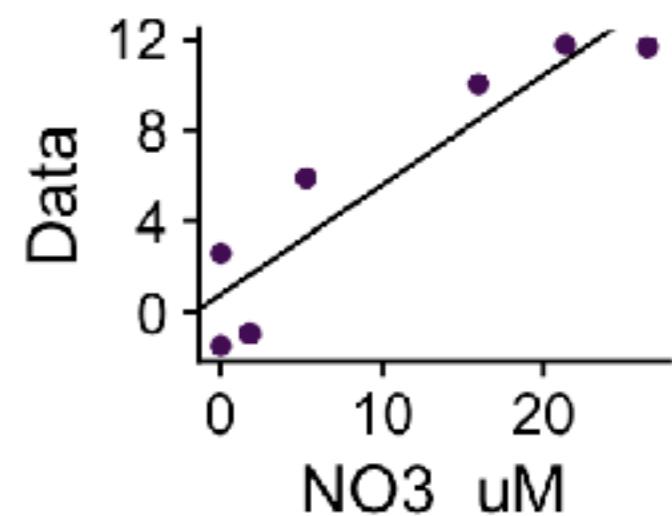
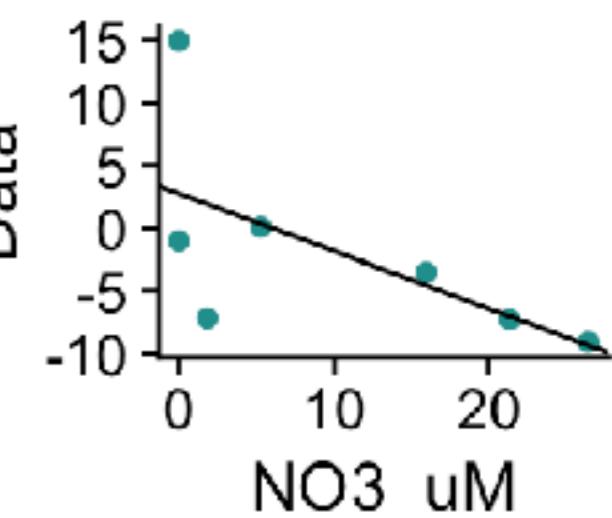
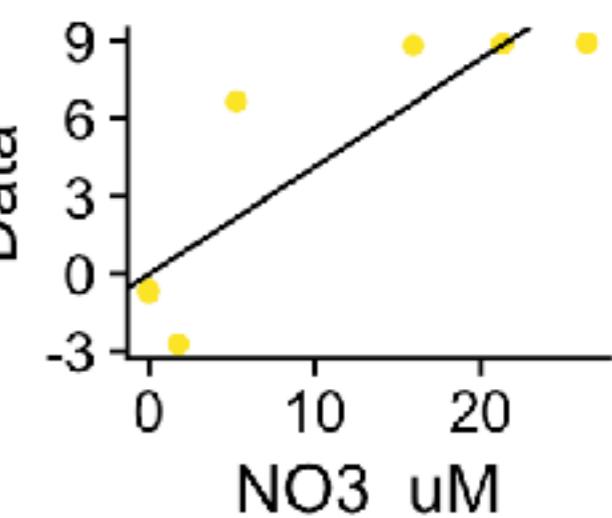
```



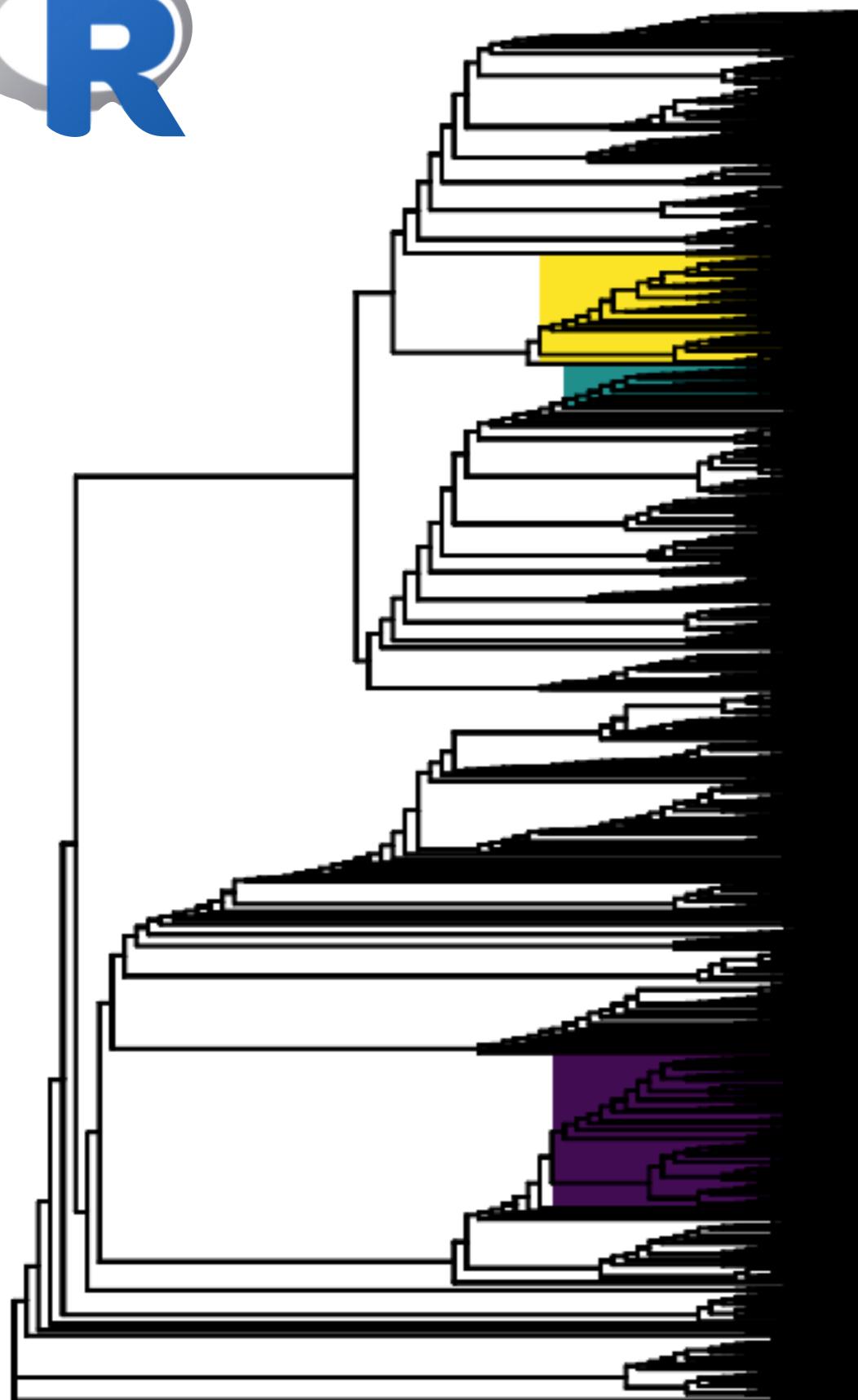




Phylogeny

**A****B****C**

R



R

