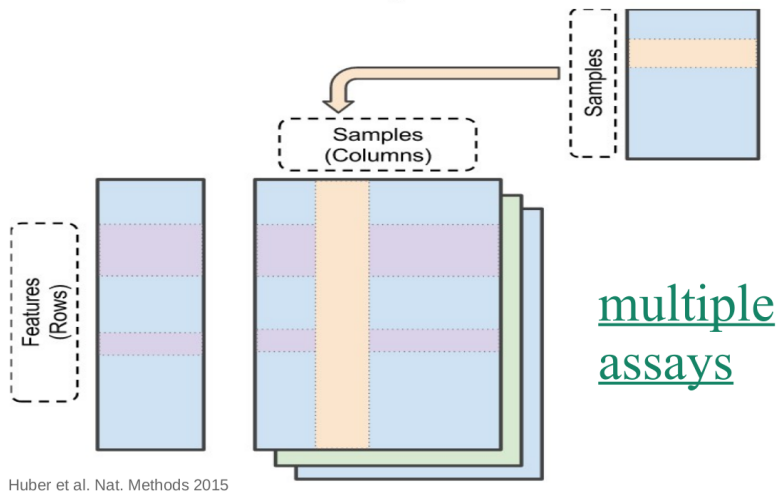# Tabular data analysis

## Leo Lahti

### Recap of Day 1

- **Data**: *SummarizedExperiment*: rowData, colData, assays
- **Methods**: subsetting, transformations



Huber et al. Nat. Methods 2015

### Day 1: Basic data wrangling

- reproducible data science workflow
- data import
- data containers
- data manipulation (subsetting, transformations)

**Today's learning goals**

Expanding multi-assay analyses:
**TreeSummarizedExperiment**

- augmenting the data (add diversities)
- data agglomeration & *alternative experiments*
- tree information: *rowTree, colTree*

**Today's program**

Morning: data wrangling

Afternoon: data visualizations

# Data enrichment

### Visualizing colData

Task: visualize the abundance of a specific microbial Species against the measurement Site
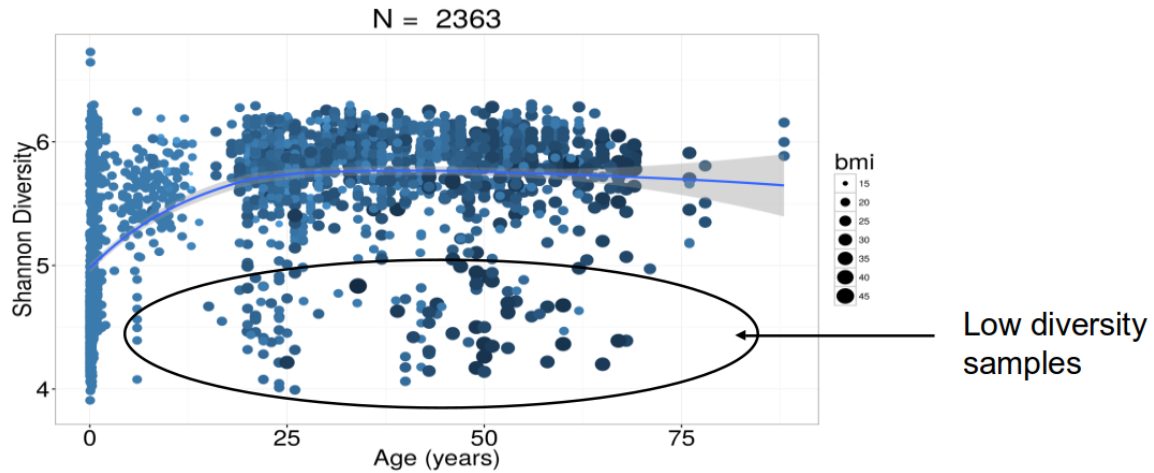
### Alpha diversity task

Use the available tools to assess and visualize alpha diversity, and augment colData

- Exercises 17.5.1-17.5.2
- Add Shannon diversity in colData
- Visualize diversity differences between sample groups

### Alpha diversity & aging

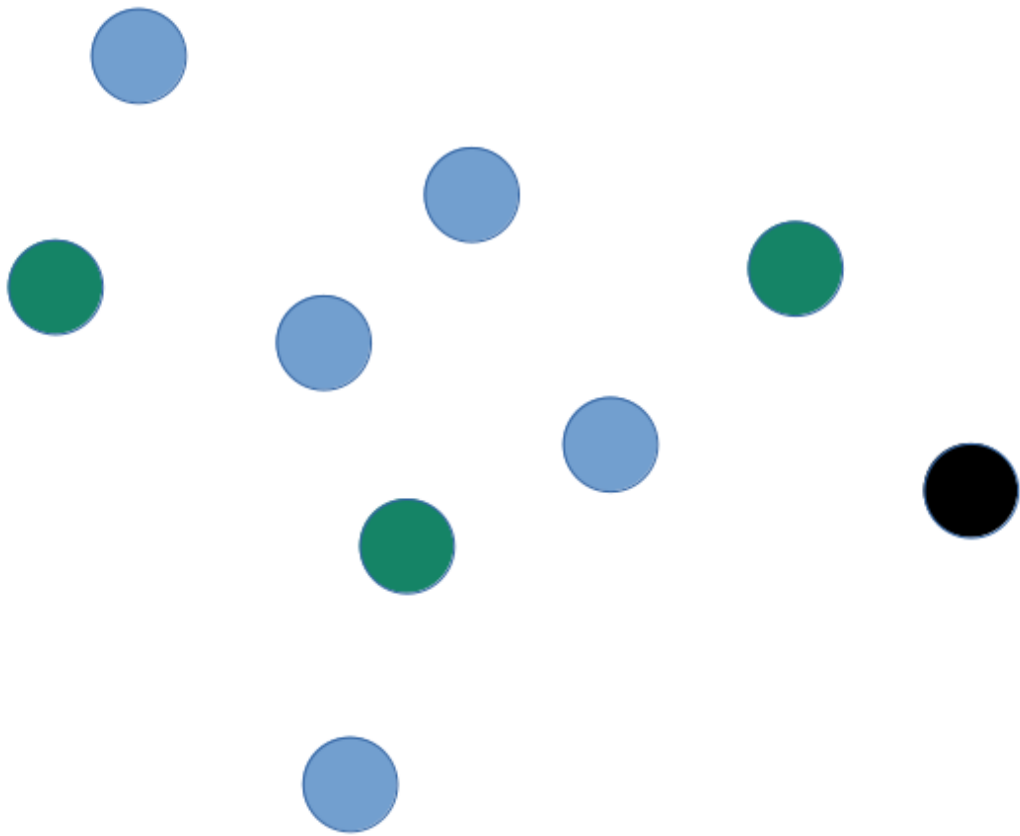Healthy & normal obese subjects.

## Alpha diversity and diet

## Associations of healthy food choices with gut microbiota profiles

Kari K Koponen, Aaro Salosensaari, Matti O Ruuskanen, Aki S Havulinna, Satu Männistö, Pekka Jousilahti, Joonatan Palmu, Rodolfo Salido, Karenina Sanders, Caitriona Brennan, Gregory C Humphrey, Jon G Sanders, Guillaume Meric, Susan Cheng, Michael Inouye, Mohit Jain, Teemu J Niiranen, Liisa M Valsta, Rob Knight, and Veikko V Salomaa

| | | |
|---|---|---|
| Combined fiber sources score | ✳ P = 0.00053 | ✳ P ≤ 0.001 |
| HFC score | ✳ P = 0.0022 | ✳ P ≤ 0.001 |
| Vegetables | P = 0.064 | ✳ P ≤ 0.001 |
| Berries | ✳ P = 0.021 | ✳ P ≤ 0.001 |
| Fruits | ✳ P = 0.014 | ✳ P ≤ 0.001 |
| Fiber-rich breads | ✳ P = 0.0027 | ✳ P ≤ 0.001 |
| Dressings and oils | P = 0.67 | ✳ P ≤ 0.001 |
| Low-fat cheeses | ✳ P = 0.015 | ✳ P = 0.02 |
| Poultry | ✳ P = 0.013 | ✳ P ≤ 0.001 |
| Red meat products (low use) | P = 0.72 | ✳ P = 0.004 |
| Juices | P = 0.86 | ✳ P = 0.023 |
| Fish | P = 0.47 | ✳ P = 0.028 |
| Nuts and seeds | P = 0.16 | P = 0.102 |

Regression coefficient in linear model for Shannon index

- −0.02, < 0.00
- 0.00, < 0.02
- 0.02, < 0.04
- 0.04, < 0.06

α

Coefficient of determination for beta diversity: 0.00% — 0.10% — 0.20% — 0.30%
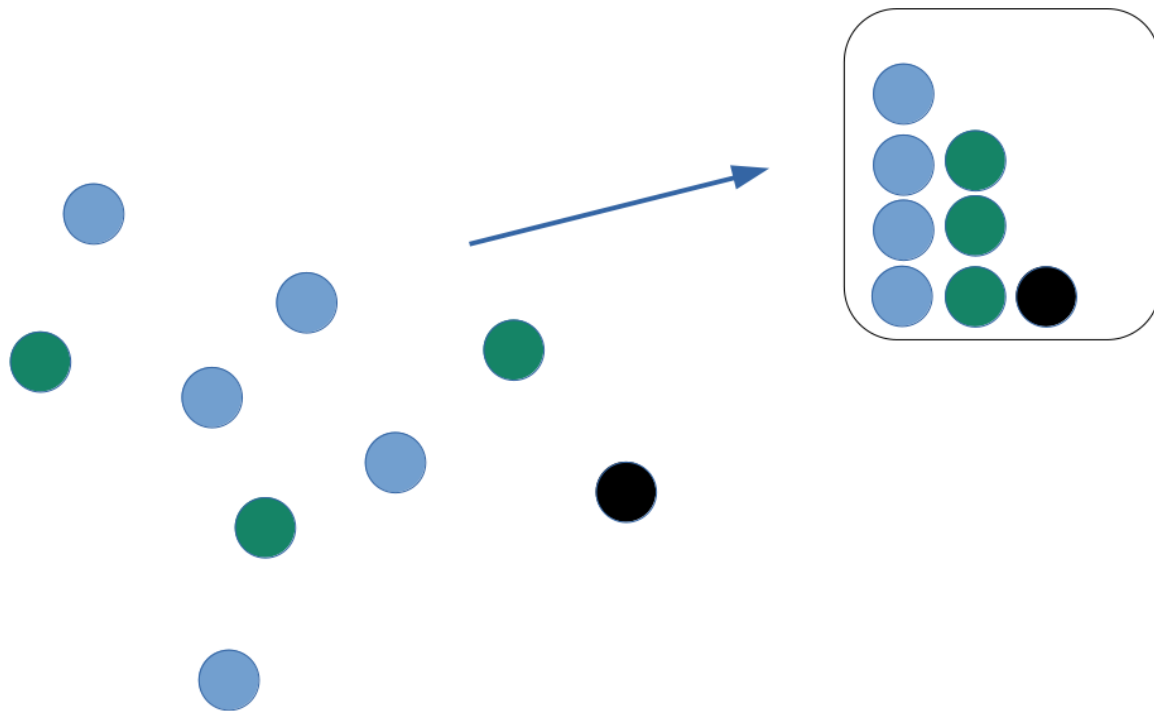
## Alpha diversity

- How many types?
- Distribution of types?
- Dominance of types?

4

**Alpha diversity**

- How many types?
- Distribution of types?
- Dominance of types?

## Alpha diversity indices

**Richness**

- number of types
- Eetimates of true richness based on finite sample sizes (Howard Sanders 1968); see e.g. Chao1

**Evenness**

- distribution of sizes (even or uneven?)
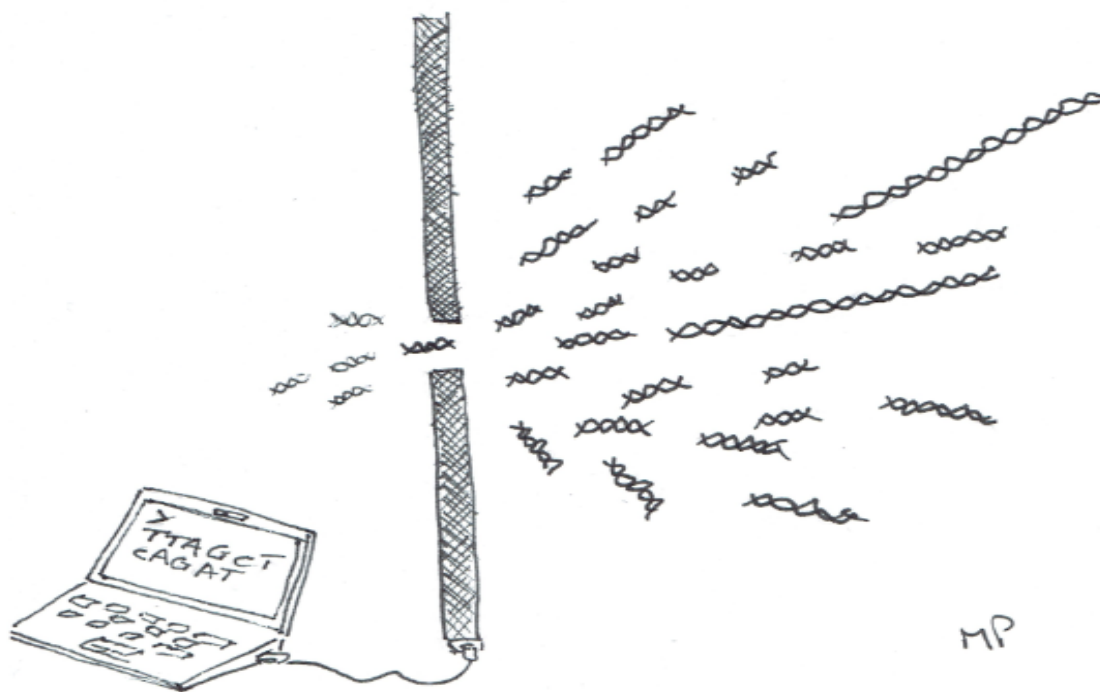
**Diversity**

- Combining richness & evenness

**Dominance**

Figure 1: https://github.com/mblstamps/stamps2019/blob/master/STAMPS2019_overview_Pop.pdf
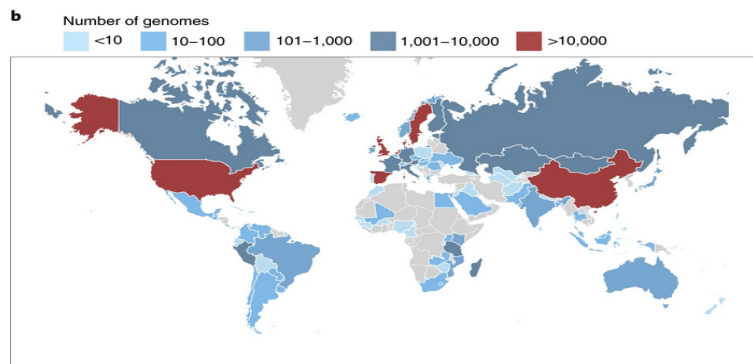
**Finite sampling**

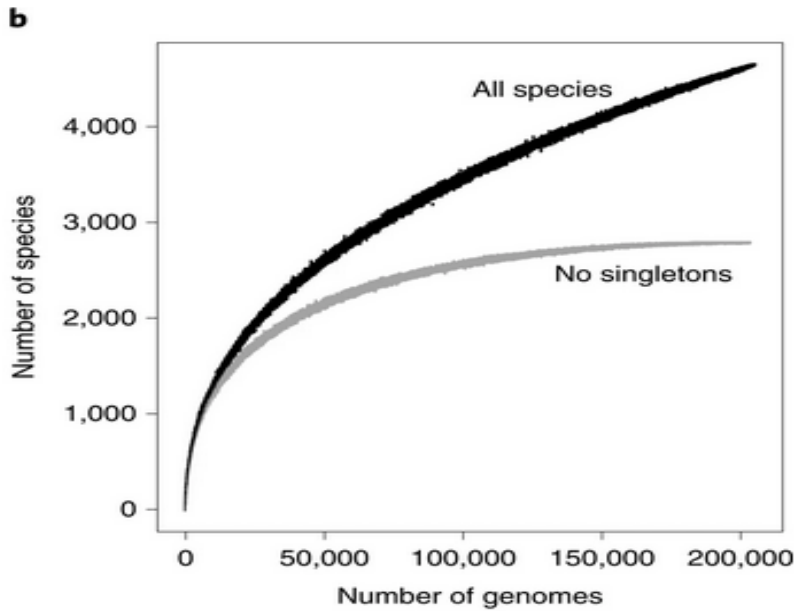## A unified catalog of 204,938 reference genomes from the human gut microbiome

Alexandre Almeida ✉, Stephen Nayfach, Miguel Boland, Francesco Strozzi, Martin Beracochea, Zhou Jason Shi, Katherine S. Pollard, Ekaterina Sakharova, Donovan H. Parks, Philip Hugenholtz, Nicola Segata, Nikos C. Kyrpides & Robert D. Finn ✉

**b**



High-quality reference genomes are required for functional characterization and taxonomic assignment of the human gut microbiota.

Unified Human Gastrointestinal Genome (UHGG):

- 4,644 gut prokaryotes (>70% lack cultured representatives)

- 204,938 nonredundant genomes

- Encode >170 million protein sequences, collated into Unified Human Gastrointestinal Protein (UHGP) catalog.

UHGP more than doubles the number of gut proteins in comparison to those present in the Integrated Gene Catalog.

- 40% of the UHGP lack functional annotations

- Intraspecies genomic variation analyses revealed a large reservoir of accessory genes and single-nucleotide variants, many of which are specific to individual human populations.

The UHGG and UHGP collections enable studies linking genotypes to phenotypes in the human gut microbiome.
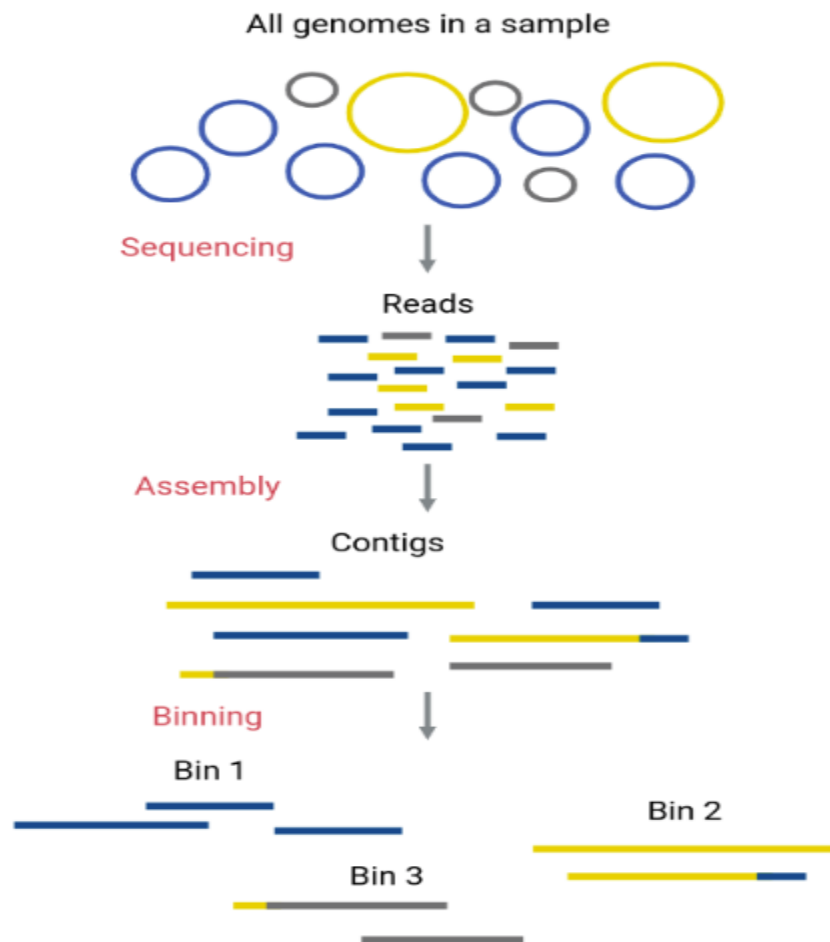
**Estimating species content**



Figure 2: Copyright © Claudia Zirion, Diego Garfias, Vanessa Arellano, Aaron Jaime, Abel Lovaco, Daniel Díaz, Abraham Avelar, Nelly Sélem https://carpentries-incubator.github.io/metagenomics-workshop/)
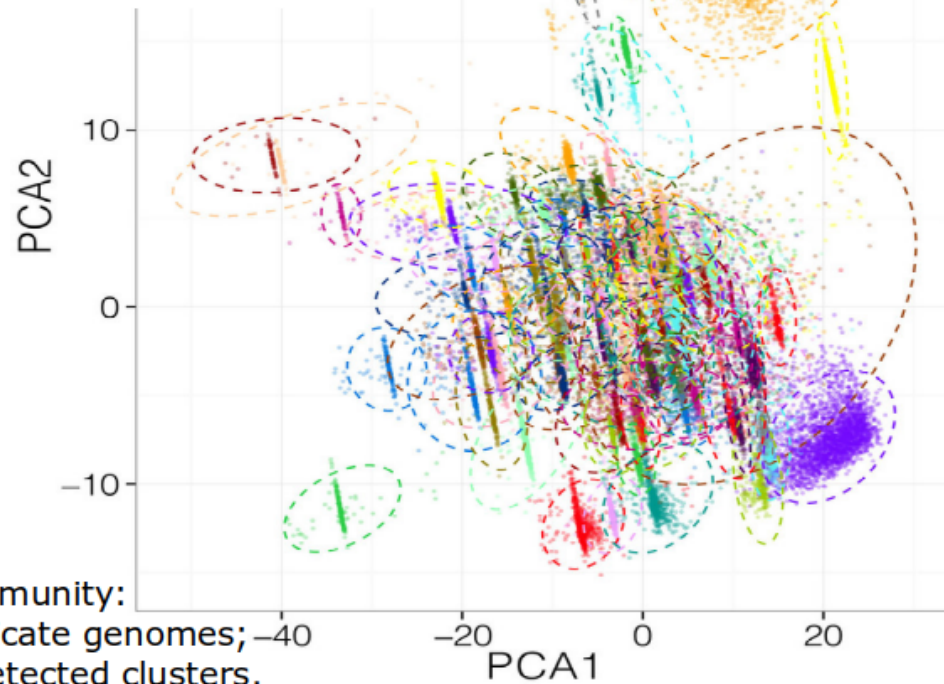
# Binning metagenomic contigs by coverage and composition

Johannes Alneberg, Brynjar Smári Bjarnason, Ino de Bruijn, Melanie Schirmer, Joshua Quick, Umer Z Ijaz, Leo Lahti, Nicholas J Loman, Anders F Andersson ✉ & Christopher Quince ✉

Mock community: colors indicate genomes; ellipses detected clusters.
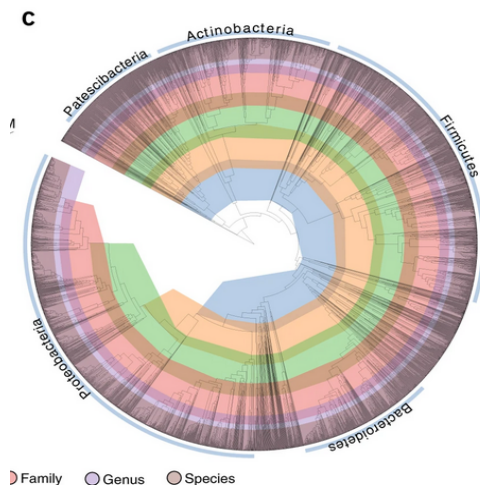
11

# A standardized bacterial taxonomy based on genome phylogeny substantially revises the tree of life

Donovan H Parks, Maria Chuvochina, David W Waite, Christian Rinke, Adam Skarshewski, Pierre-Alain Chaumeil & Philip Hugenholtz ✉

**32k** Accesses | **728** Citations | **520** Altmetric | Metrics



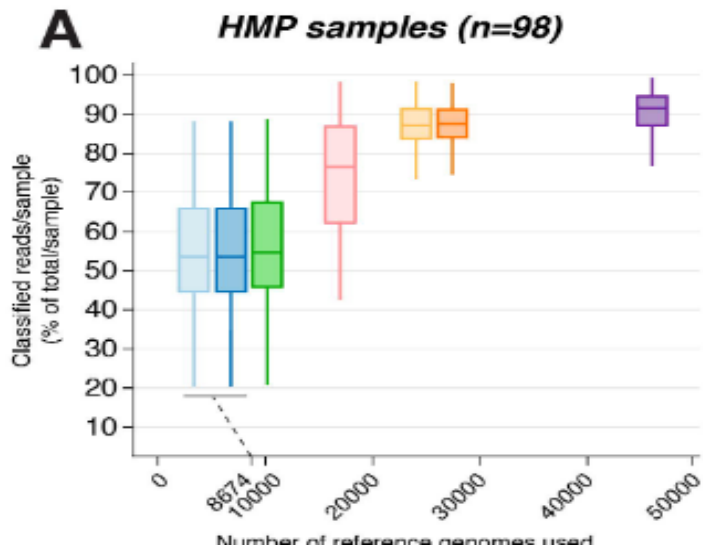# Correcting index databases improves metagenomic studies

🆔 Guillaume Méric, 🆔 Ryan R. Wick, Stephen C. Watts, 🆔 Kathryn E. Holt, 🆔 Michael Inouye

## Common alpha diversity indices

### Phylogenetically neutral diversities:

- Richness (observed, Chao1, ACE)
- Evenness (Pielou's evenness)
- Diversity (inverse Simpson, Shannon)

### Phylogeny-aware diversities:

- Faith diversity index

13

**Phylogenetic diversity indices**



# A guide to phylogenetic metrics for conservation, community ecology and macroecology

Caroline M. Tucker ✉, Marc W. Cadotte, Silvia B. Carvalho, T. Jonathan Davies, Simon Ferrier
, Susanne A. Fritz, Rich Grenyer, Matthew R. Helmus, Lanna S. Jin ... See all authors ∨

## Inverse Simpson

$$DI = \frac{N(N-1)}{\sum n(n-1)}$$

**KEY** →

$N$ = Total number of individuals collected

$n_i$ = Number of individuals of a species

$DI$ = Simpson Diversity Index

How likely it is to pick two members of the same species at random?

**Inverse Simpson**

**Beware** the variants:

- Simpson $(\lambda)$
- reciprocal Simpson $(1 - \lambda)$
- inverse Simpson $\left(\frac{1}{\lambda}\right)$

**Shannon diversity**

Shannon Index:

$$H' = -\sum_{i=1}^{S} p_i \ln p_i$$

True Richness:

$$\exp(\mathrm{H})$$

*True diversity, or the effective number of types, refers to the number of equally abundant types needed for the average proportional abundance of the types to equal what is observed in the dataset of interest.*

**Evenness**

H / ln(S)

- H: Shannon diversity
- S: Species richness

**Hill's Diversity as a unifying concept**

$$^{q}D = \left(\sum_{i}^{R} p_i^q\right)^{\frac{1}{1-q}} \tag{1}$$

**Hill's alpha diversities**

R: richness (number of distinct types)

pi: proportion of type I

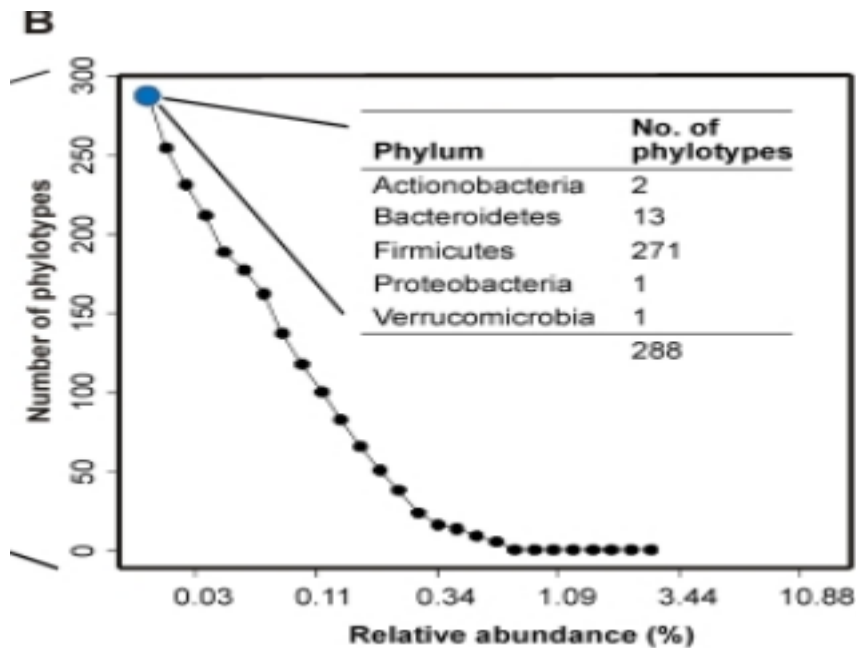Order of diversity:

- q = 0 : Species Richness
- q = 1 : Shannon diversity
- q = 2 : (Inverse) Simpson diversity
- q  1 : Renyi entropy

## Hill's Diversity as a unifying concept



## Hill's alpha diversities

- Richness
- inverse Simpson
- Shannon

# Data wrangling

## Basic data operations

- Transform

17

- Subset
- Merge
- **Aggregate**
- Split

## Subsetting

Load example data set:

```
library(mia)
```

Loading required package: SummarizedExperiment

Loading required package: MatrixGenerics

Loading required package: matrixStats


Attaching package: 'MatrixGenerics'

The following objects are masked from 'package:matrixStats':

    colAlls, colAnyNAs, colAnys, colAvgsPerRowSet, colCollapse,
    colCounts, colCummaxs, colCummins, colCumprods, colCumsums,
    colDiffs, colIQRDiffs, colIQRs, colLogSumExps, colMadDiffs,
    colMads, colMaxs, colMeans2, colMedians, colMins, colOrderStats,
    colProds, colQuantiles, colRanges, colRanks, colSdDiffs, colSds,
    colSums2, colTabulates, colVarDiffs, colVars, colWeightedMads,
    colWeightedMeans, colWeightedMedians, colWeightedSds,
    colWeightedVars, rowAlls, rowAnyNAs, rowAnys, rowAvgsPerColSet,
    rowCollapse, rowCounts, rowCummaxs, rowCummins, rowCumprods,
    rowCumsums, rowDiffs, rowIQRDiffs, rowIQRs, rowLogSumExps,
    rowMadDiffs, rowMads, rowMaxs, rowMeans2, rowMedians, rowMins,
    rowOrderStats, rowProds, rowQuantiles, rowRanges, rowRanks,
    rowSdDiffs, rowSds, rowSums2, rowTabulates, rowVarDiffs, rowVars,
    rowWeightedMads, rowWeightedMeans, rowWeightedMedians,
    rowWeightedSds, rowWeightedVars

```
Loading required package: GenomicRanges

Loading required package: stats4

Loading required package: BiocGenerics


Attaching package: 'BiocGenerics'

The following objects are masked from 'package:stats':

    IQR, mad, sd, var, xtabs

The following objects are masked from 'package:base':

    anyDuplicated, aperm, append, as.data.frame, basename, cbind,
    colnames, dirname, do.call, duplicated, eval, evalq, Filter, Find,
    get, grep, grepl, intersect, is.unsorted, lapply, Map, mapply,
    match, mget, order, paste, pmax, pmax.int, pmin, pmin.int,
    Position, rank, rbind, Reduce, rownames, sapply, setdiff, sort,
    table, tapply, union, unique, unsplit, which.max, which.min

Loading required package: S4Vectors


Attaching package: 'S4Vectors'

The following objects are masked from 'package:base':

    expand.grid, I, unname

Loading required package: IRanges

Loading required package: GenomeInfoDb

Loading required package: Biobase
```

```
Welcome to Bioconductor

    Vignettes contain introductory material; view with
    'browseVignettes()'. To cite Bioconductor, see
    'citation("Biobase")', and for packages 'citation("pkgname")'.


Attaching package: 'Biobase'

The following object is masked from 'package:MatrixGenerics':

    rowMedians

The following objects are masked from 'package:matrixStats':

    anyMissing, rowMedians

Loading required package: SingleCellExperiment

Loading required package: TreeSummarizedExperiment

Loading required package: Biostrings

Loading required package: XVector


Attaching package: 'Biostrings'

The following object is masked from 'package:base':

    strsplit

Loading required package: MultiAssayExperiment
```

```r
data(GlobalPatterns)
tse <- GlobalPatterns
```

Check dimension:

```
dim(tse)
```

```
[1] 19216    26
```

Check dimension for a subset:

```
dim(tse[1:10, 1:3])
```

```
[1] 10  3
```

## Transformations

- Presence/absence
- Compositional (percentages)
- $Log_{10}$
- CLR and other *Aitchison* transformations
- Phylogenetic transformations (e.g. philr)
- Custom transformations

## Transformations

Task: Alternative assays

- visualize transformed data; histograms, boxplots
- compare different transformations (scatterplot?)

## Agglomeration

- taxonomic units
- TreeSE objects

## Agglomeration

Agglomerate microbiota data to higher taxonomic levels:

- chapter 6.3
- agglomerateByRank
- compare diversity or prevalent features between levels

## Alternative experiments

Alternative assays vs. alternative experiments?

- Store agglomerated data: *altExp*
- Do all levels at once: *splitByRanks*

## Splits

Splitting by:

- taxonomic units
- sample or feature groups

## Taxonomic ranks & *altExp*

The alternative experiments (*altExp*) mechanism allows us to include multiple abundance tables at different taxonomic levels.
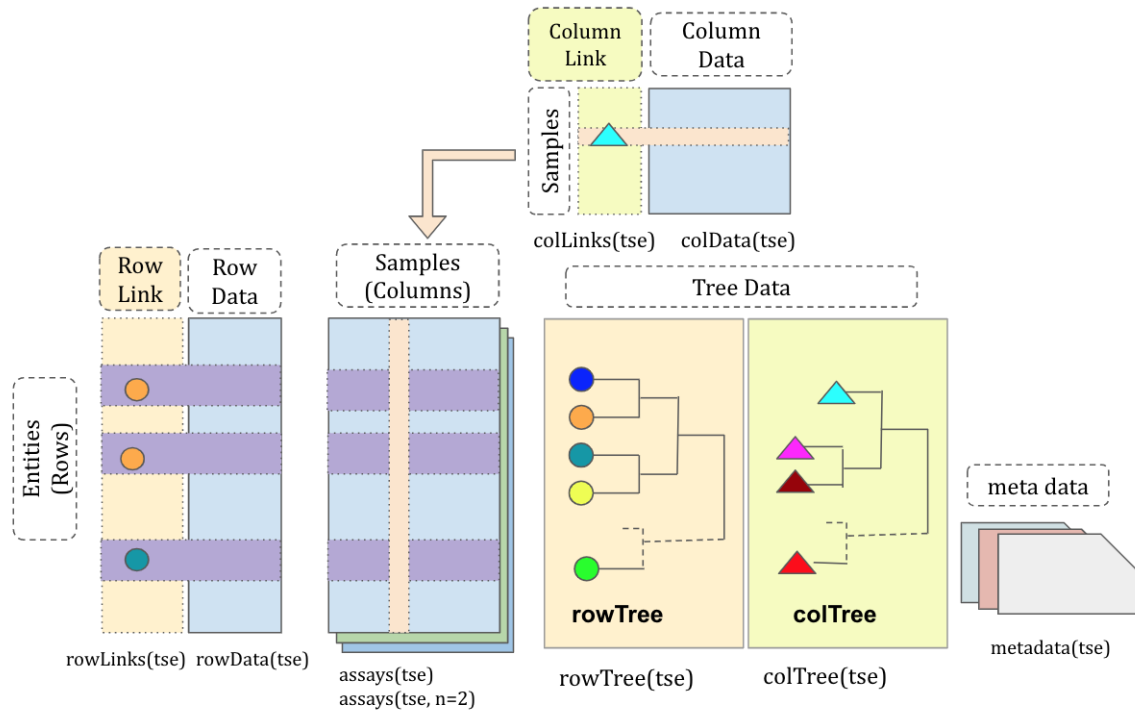
| Option | Rows (features) | Cols (samples) | Recommendation |
|---:|---|---|---|
| assays | match | match | Data transformations |
| **altExp** | free | match | Alternative experiments |
| MultiAssay | free | free (mapping) | Multi-omic experiments |

## Alternative experiments and assays?

- Pick clr assay from Genus-level data table?
- Compare Shannon diversity from Genus and Species levels?

## TreeSummarizedExperiment

Huang et al. F1000, 2021

colLinks(tse)  colData(tse)

Tree Data

rowTree  colTree

meta data

rowLinks(tse)  rowData(tse)

assays(tse)
assays(tse, n=2)

rowTree(tse)  colTree(tse)

metadata(tse)

# Visualization

## Ordination

- Visualize example data with PCoA using Bray-Curtis dissimilarity
- Visualize example data with PCA using Aitchison distance (CLR + Euclid)

## Heatmaps

- Visualize abundance variation for selected taxa on a heatmap

## Trees

- Visualize phylogenetic tree using the examples