

Microbiome data science workflow

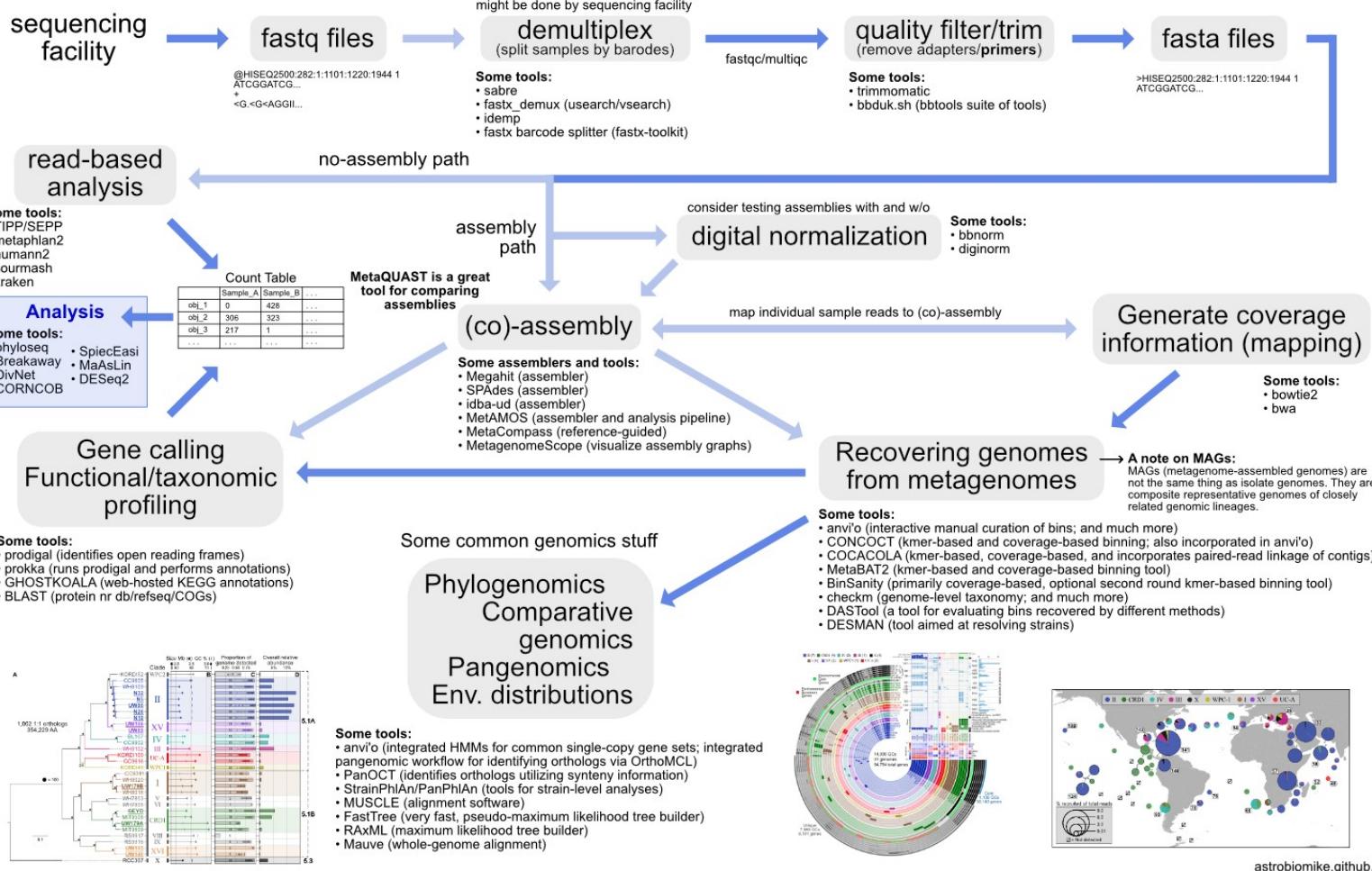
Leo Lahti, CSC course, Nov 28 – Dec 2, 2022



Overview of generic* metagenomics workflow

*This is generic; specific workflows can vary on the order of steps here and how they are done.

When working with your own data you should never follow any pipeline blindly. There can be critical differences based on your data.



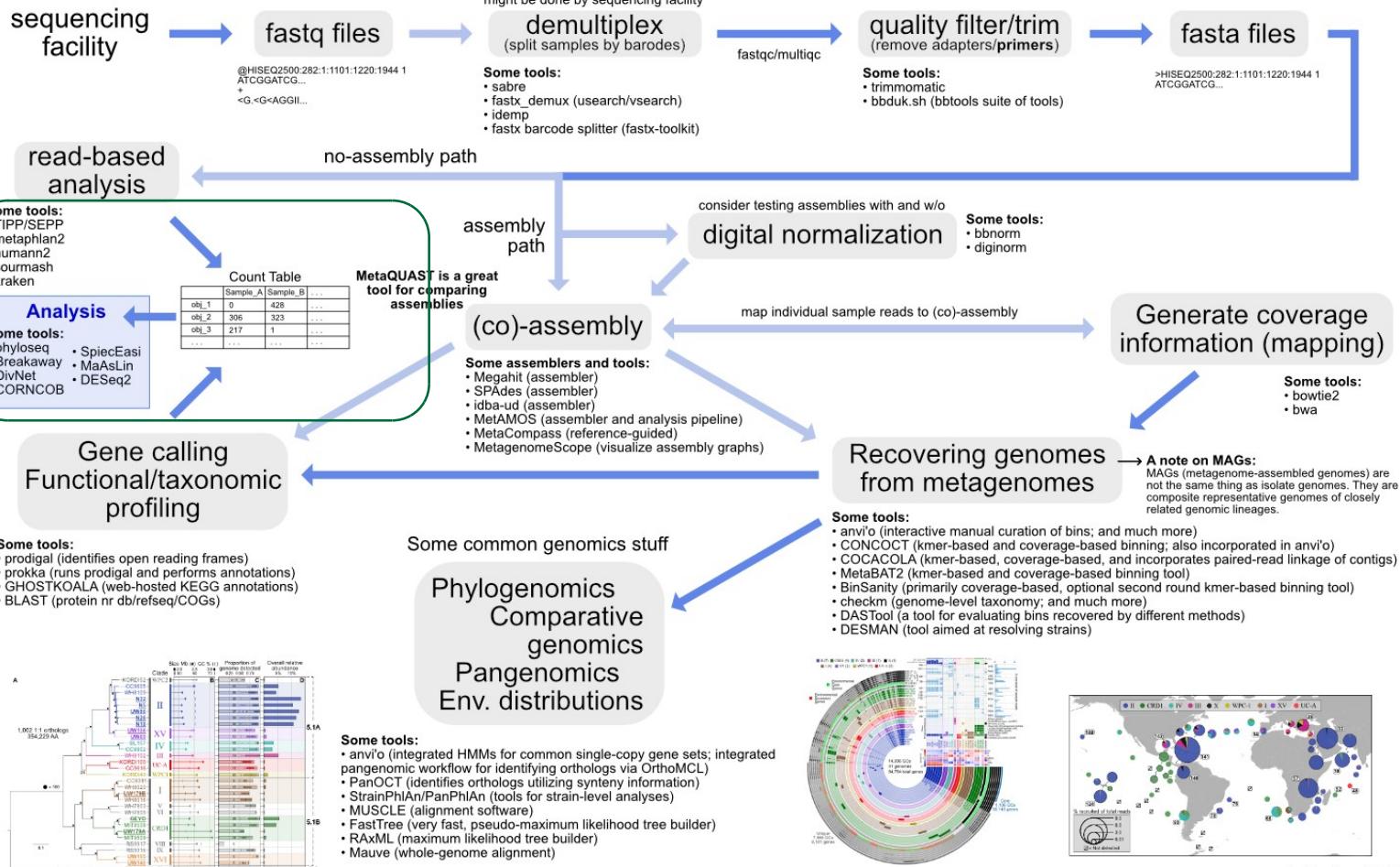
Happy Belly Bioinformatics
JOSE 10.21105/jose.00053

AstroBioMike
Orcid: 0000-0001-7750-9145

Lee, (2019). Happy Belly Bioinformatics: an open-source resource dedicated to helping biologists utilize bioinformatics. Journal of Open Source Education, 4(41), 53, <https://doi.org/10.21105/jose.00053>

Overview of generic* metagenomics workflow

*This is generic; specific workflows can vary on the order of steps here and how they are done.



Our primary focus is on downstream analysis of taxonomic profiles of the microbial community, obtained for instance with:

- 16S rRNA amplicon sequencing
- (shallow) shotgun sequencing
- deep metagenome sequencing
- phylogenetic microarrays

Happy Belly Bioinformatics

JOSE 10.21105/jose.00053

AstroBioMike
Orcid: 0000-0001-7750-9145

Lee, (2019). Happy Belly Bioinformatics: an open-source resource dedicated to helping biologists utilize bioinformatics. Journal of Open Source Education, 4(41), 53, <https://doi.org/10.21105/jose.00053>

Common study designs

Case-control & Intervention
targeted experimental testing

Cross-sectional
population (cohort) studies

Prospective
long-term follow-ups

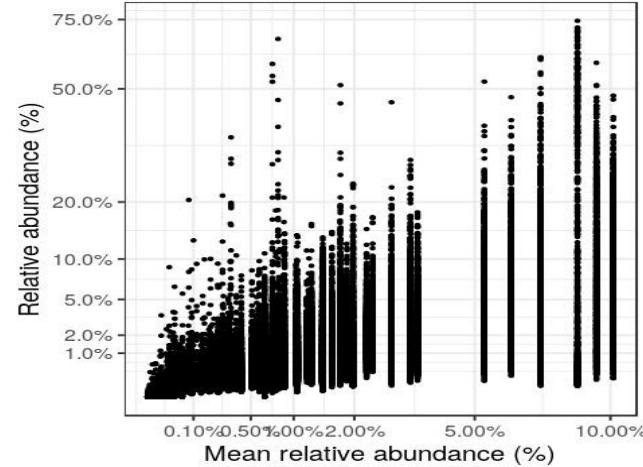
Longitudinal
ecosystem dynamics

Microbiome data properties

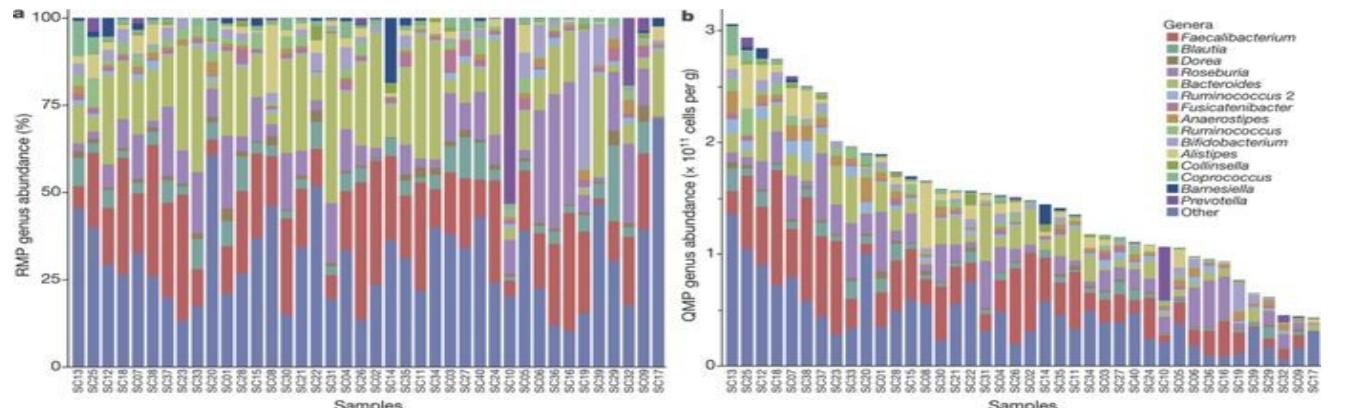
- Sparse
- Non-Gaussian
- Overdispersed
- Compositional
- Complex
- Stochastic
- Hierarchical

Heteroschedasticity

Data: HITChip Atlas (Lahti et al.
Nature Communications 2014)



Compositionality



open microbiome data science frameworks

mothur

Download Wiki Forum Blog Riggmanes Facebook

Welcome to the website for the mothur project, initiated by Dr. Patrick Schloss and his research group in the Department of Microbiology & Immunology at the University of Michigan. This project seeks to develop a single piece of open-source, expandable software to fit the bioinformatics needs of the microbial community. In February 2009 we released the first version of mothur, which had accelerated versions of the most common bioinformatics programs. mothur has gone on to become one of the most cited bioinformatics tools for analyzing 16S rRNA gene experiments. They now have a mailing forum and team how you can use mothur to process data from samples from Bangladesh, India, and Bhutan (MothurSeeds). If you would like to contribute code to the project feel free to download the source code and make your own improvements. Alternatively, if you have an idea of a need, but lack the programming expertise, let us know through the forums and we will add it to the queue of features we would like to add.

Subscribe to the mothur mailing list
email address _____
Subscribe

Department of Microbiology & Immunology
The University of Michigan Medical School
The University of Michigan
This site is maintained by Pat Schloss
© 2008-2010

QIIME 2™ is a next-generation microbiome bioinformatics platform that is extensible, free, open source, and community developed.

Code of Conduct » Citing QIIME 2 » Learn more »

Automatically track your analyses with decentralized data provenance — no more guesswork on what commands were run!

Interactively explore your data with beautiful visualizations that provide new perspectives.

Easily share results with your team, even those members without QIIME 2 installed.

Plugin-based system — your favorite microbiome methods all in one place.



[PeerJ >](#)

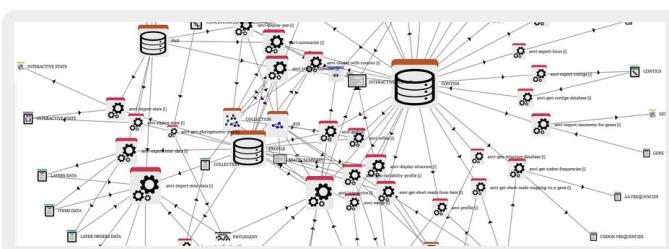
Anvi'o: an advanced analysis and visualization platform for 'omics data

Research article Bioinformatics Biotechnology Computational Biology Genomics Microbiology

A. Murat Eren^{✉ 1,2}, Özcan C. Esen¹, Christopher Quince³, Joseph H. Vineis¹, Hilary G. Morrison¹, Mitchell L. Sogin¹, Tom O. Delmont¹

Published October 8, 2015

Anvi'o in a nutshell



Anvi'o is an open-source, community-driven analysis and visualization platform for 'omics data.

Workflow for community analysis

Data import

- Container preparation
- Sanity check (components)

Data analysis

- Alpha & beta diversity
- Differential abundance
- Etc.

Reporting

- Summaries
- Reproducible reporting

F1000Research

F1000Research 2016, 5:1492 Last updated: 02 AUG 2016

Check for updates

RESEARCH ARTICLE

RESEARCH ARTICLE

REVISED Bioconductor Workflow for Microbiome Data Analysis: from raw reads to community analyses [version 2; peer review: 3 approved]

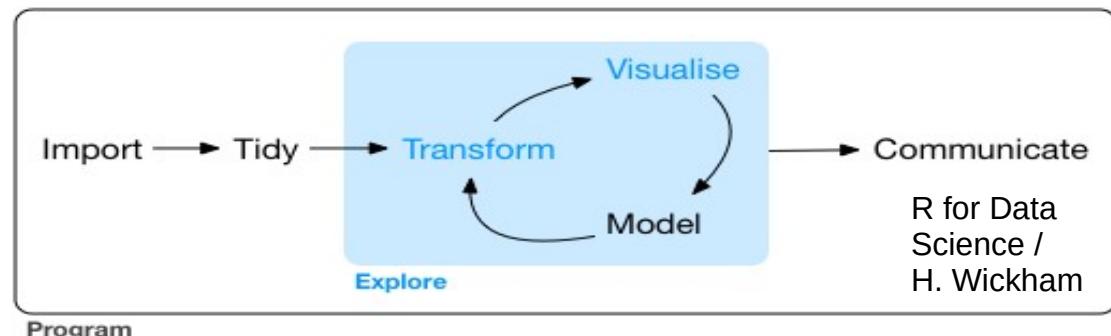
Ben J. Callahan¹, Kris Sankaran¹, Julia A. Fukuyama¹, Paul J. McMurdie², Susan P. Holmes¹

¹Stanford University, Stanford, CA, 94305, USA

²Whole Biome Inc., San Francisco, CA, 94107, USA

¹Statistics Department, Stanford University, Stanford, CA, 94305, USA

²Whole Biome Inc., San Francisco, CA, 94107, USA



Workflow for community analysis

Data import

- Container preparation
- Sanity check (components)

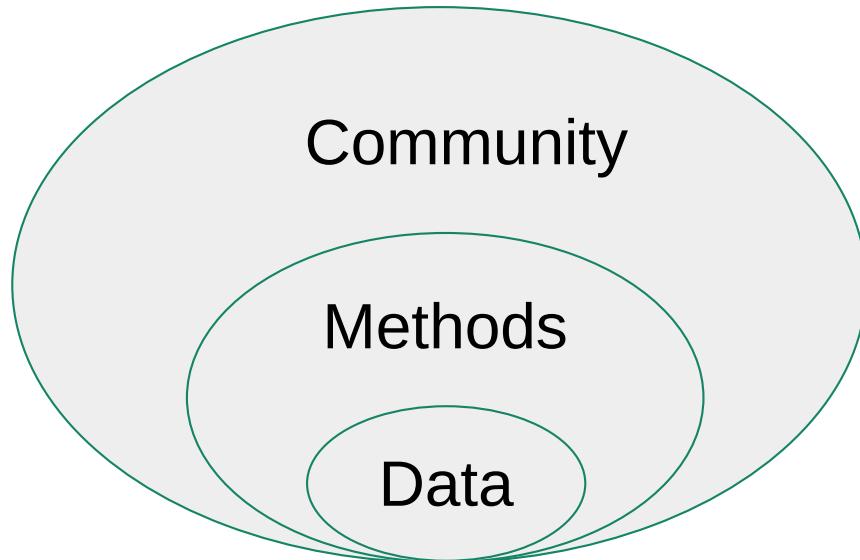
Data analysis

- Alpha & beta diversity
- Differential abundance
- Etc.

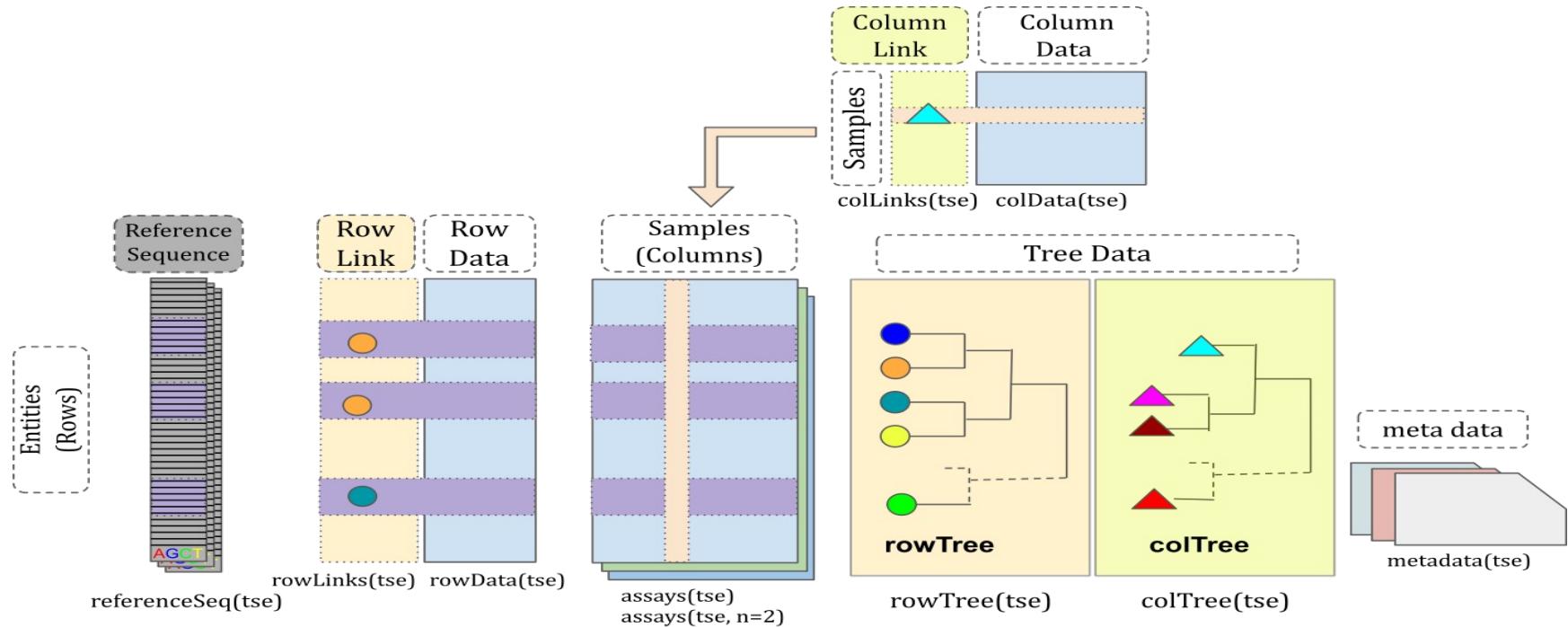
Reporting

- Summaries
- Reproducible reporting

Standardized data containers are central for the R/Bioc ecosystem



TreeSummarizedExperiment data container



Data analysis & visualization

Data import

- Container preparation
- Sanity check (components)

```
plotColData(tse,
            "observed",
            "SampleType",
            colour_by = "Final_Barcodes") +
  theme(axis.text.x = element_text(angle=45,hjust=1)) +
  ylab(expression(Richness[Observed]))
```

Data analysis

- Alpha & beta diversity
- Differential abundance
- Etc.

Reporting

- Summaries
- Reproducible reporting

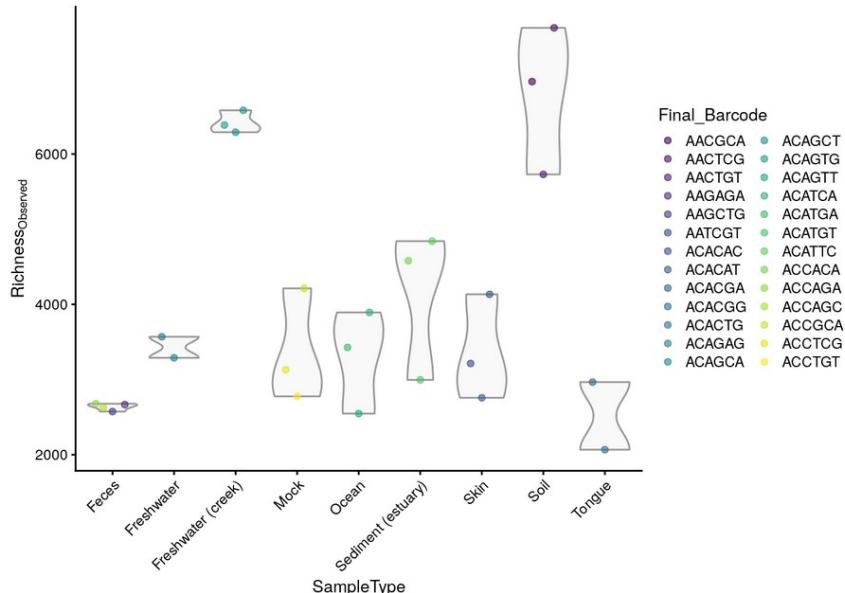


Figure 7.1: Shannon diversity estimates plotted grouped by sample type with colour-labeled barcode.

Typical tasks in community analysis

Standard:

- Community diversity (alpha diversity)
- Community similarity (beta diversity)
- Differential abundance analysis

Frequent:

- Community typing
- Co-occurrence networks
- Phylogenetic balances

```
plotRowTree(x[rowData(x)$Phylum %in% top_phyla_mean, ],  
            edge_colour_by = "Phylum",  
            tip_colour_by = "prevalence",  
            node_colour_by = "prevalence")
```

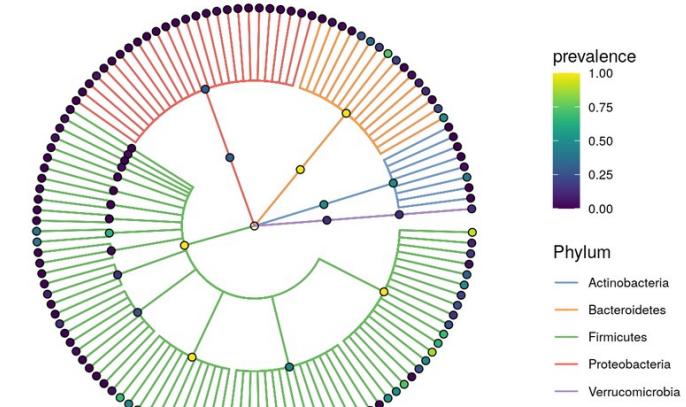
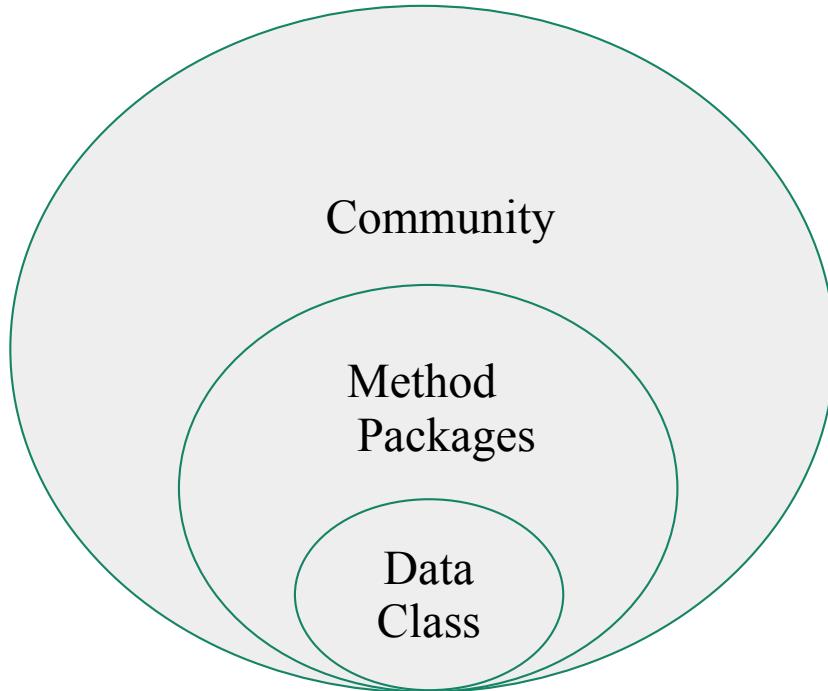


Figure 5.2: Prevalence of top phyla as judged by mean abundance

Reduce overlapping efforts, improve interoperability, ensure sustainability.



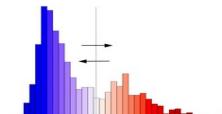
Data packages

ExperimentHub

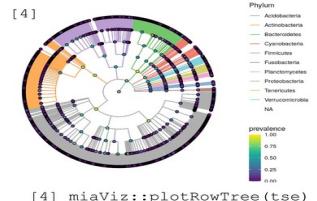
platforms all rank 76 / 1974 posts 2 / 1 / 2e+01 / 1 in Bioc 4 years
build ok updated before release dependencies 72

DOI: [10.18129/B9.bioc.ExperimentHub](https://doi.org/10.18129/B9.bioc.ExperimentHub) [f](#) [t](#)

mia – microbiome analysis
getDiversity(x)
calculateDMM(x)



miaViz - Visualization



Package ecosystem

Workflow for community analysis

Data import

- Container preparation
- Sanity check (components)

Data analysis

- Exploration
- Analysis & modeling

Reporting

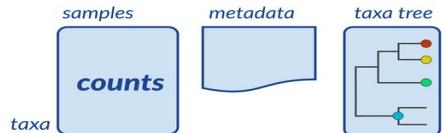
- Summaries
- Reproducible reporting

Data containers support collaborative development of analysis methods & workflows

Import Data

This workflow starts with either raw data directly from relative abundance estimation or taxonomic classification
OR pre-existing data objects from widely used software.

RAW DATA

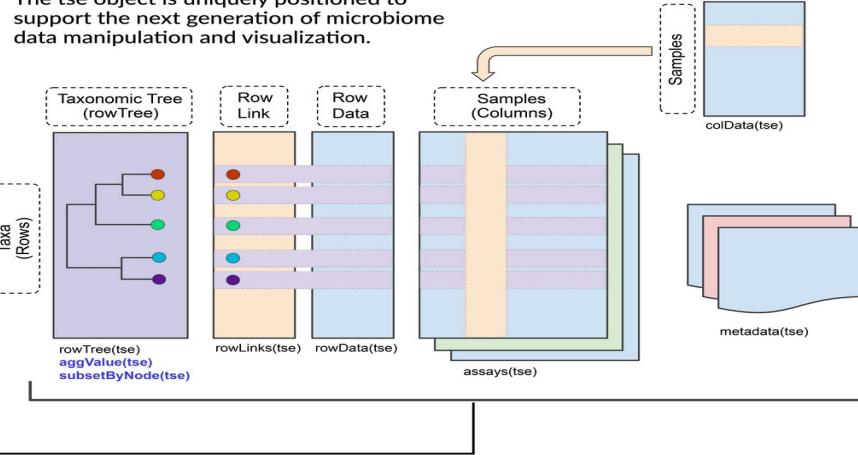


EXISTING DATA



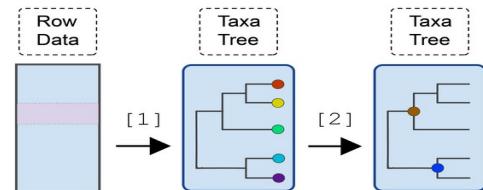
The TreeSE object

The tse object is uniquely positioned to support the next generation of microbiome data manipulation and visualization.

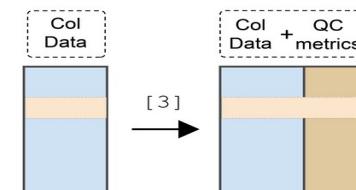


The mia Pipeline

Accessing Taxonomic Info.



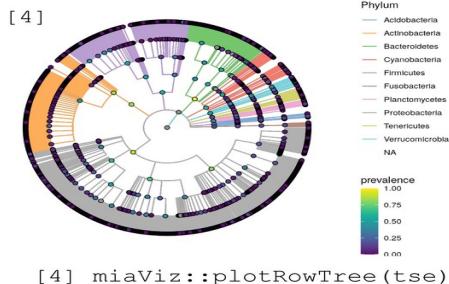
Quality Control



[1] mia::addTaxonomyTree(tse)
[2] TreeSE::aggValue(tse)

[3] scatter::addPerCellQC(tse)

Visualizing with miaViz



Orchestrating Microbiome Analysis with R and Bioconductor – online book: *beta version*



microbiome.github.io

Data

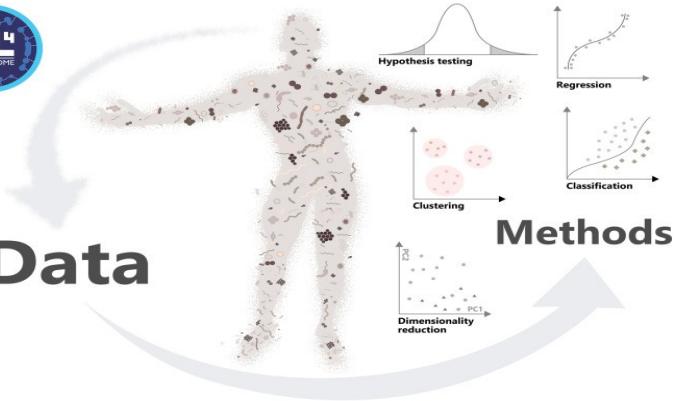


Figure source: Moreno-Indias et al. (2021) Statistical and Machine Learning Techniques in Human Microbiome Studies: Contemporary Challenges and Solutions. *Frontiers in Microbiology* 12:11.

Microbiome Analysis

Authors: Leo Lahti [aut], Sudarshan Shetty [aut], Felix GM Ernst [aut, cre]

Version: 0.98.0003

Modified: 2020-12-06

Compiled: 2020-12-13

Environment: R version 4.0.0 (2020-04-24), Bioconductor 3.11

License: CC BY-NC-SA 3.0 US

Copyright:

Source: <https://github.com/microbiome/MiaBook>

Preface

This website is a book on microbiome analysis in the Bioconductor universe and is showing common principles and workflows of performing microbiome analysis.

The book was borne out of necessity, while updating tools for microbiome analysis to work with common classes of the Bioconductor project handling count data of various sorts. It is heavily influenced by similar resources, such as the [Orchestrating Single-Cell Analysis with Bioconductor](#) book, [phyloseq tutorials](#) and [microbiome tutorials](#).

We focus on microbiome analysis tools, new, updated and established methods. In the *Introduction* section, we show how to work with the key data infrastructure `TreeSummarizedExperiment` and related classes, how this framework relates to other infrastructure and how to load microbiome analysis data to work with in the context of this framework.

The second section, *Focus Topics*, is all about the steps for analyzing microbiome data, beginning with the most common steps and progressing to more specialized methods in subsequent sections.

The third section, *Appendix*, contains the rest of things we didn't find another place for, yet.

Acknowledgments

Course organizers:

- Finnish IT Center for Science (CSC)
- Department of Computing, University of Turku, Finland

Material preparation supported by:



Turun yliopisto
University of Turku



BIOCITY TURKU
BIOCITY TURKU
BIOCITY TURKU



SUOMEN AKATEMIA
FINLANDS AKADEMI • ACADEMY OF FINLAND



MICROBIOME

