

Metabarcoding processing pipelines and considerations

Soledad Benitez Ponce (benitezponce.1@osu.edu)

March 14, 2024

Received: 10 February 2023

Revised: 5 June 2023

Accepted: 6 July 2023

DOI: 10.1111/1755-0998.13847

FROM THE COVER

MOLECULAR ECOLOGY
RESOURCES WILEY

A pile of pipelines: An overview of the bioinformatics software for metabarcoding data analyses

Ali Hakimzadeh¹  | Alejandro Abdala Asbun² | Davide Albanese³ | Maria Bernard^{4,5}  |
Dominik Buchner⁶  | Benjamin Callahan⁷ | J. Gregory Caporaso⁸ | Emily Curd⁹ |
Christophe Djemiel¹⁰  | Mikael Brandström Durling¹¹  | Vasco Elbrecht⁶  |
Zachary Gold¹² | Hyun S. Gweon^{13,14} | Mehrdad Hajibabaei¹⁵  | Falk Hildebrand^{16,17} |
Vladimir Mikryukov¹ | Eric Normandeau¹⁸ | Ezgi Özkurt^{16,17} | Jonathan M. Palmer¹⁹ |
Géraldine Pascal²⁰  | Teresita M. Porter¹⁵ | Daniel Straub²¹ | Martti Vasar¹  |
Tomáš Větrovský²² | Haris Zafeiropoulos²³ | Sten Anslan^{1,24} 

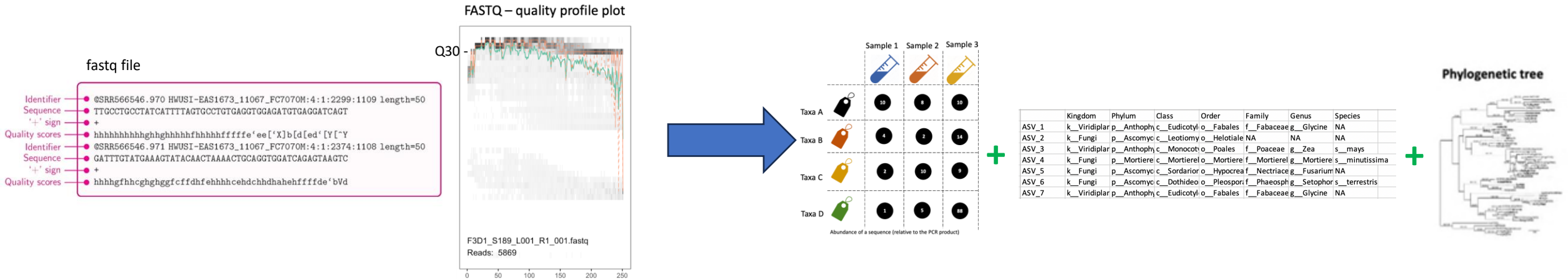
Metabarcoding processing and analysis pipelines

1. From raw-reads to OTU-table

2. Statistical analysis

Metabarcoding processing and analysis pipelines

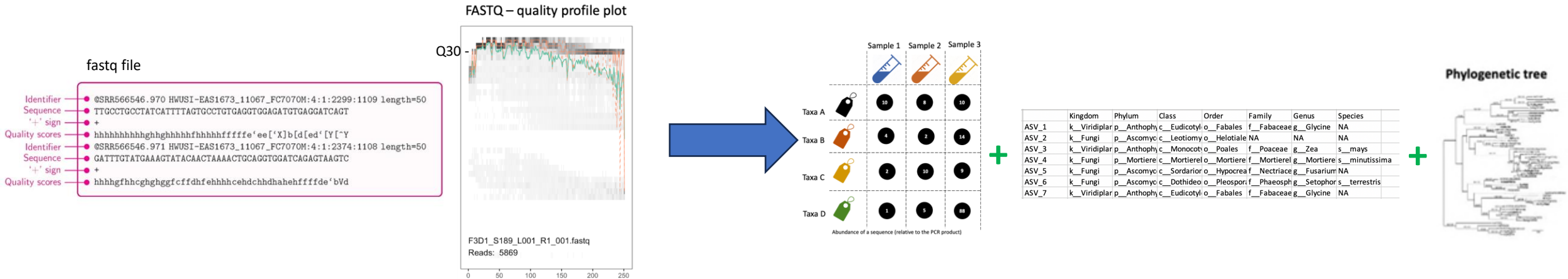
1. From raw-reads to OTU-table



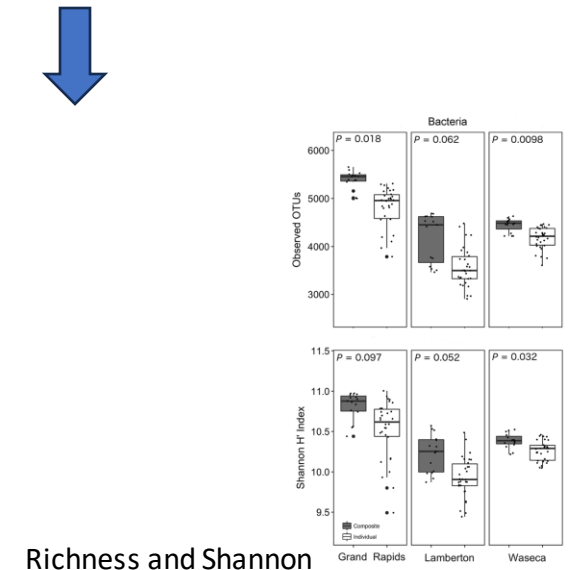
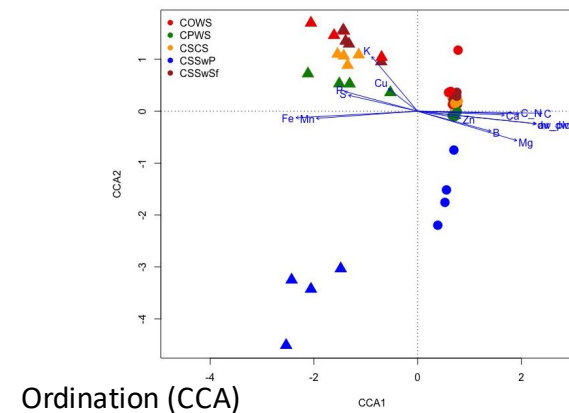
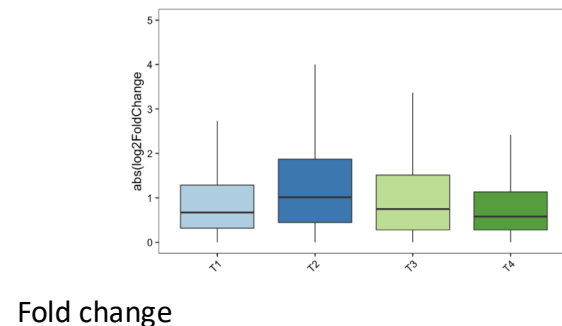
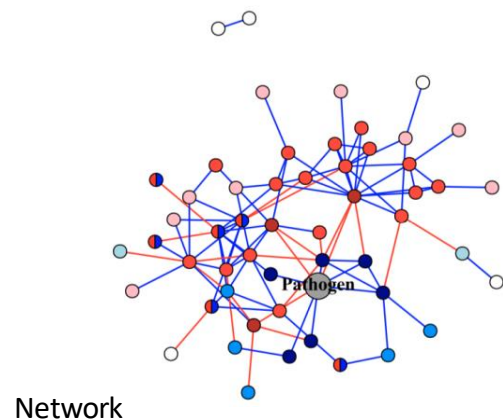
2. Statistical analysis

Metabarcoding processing and analysis pipelines

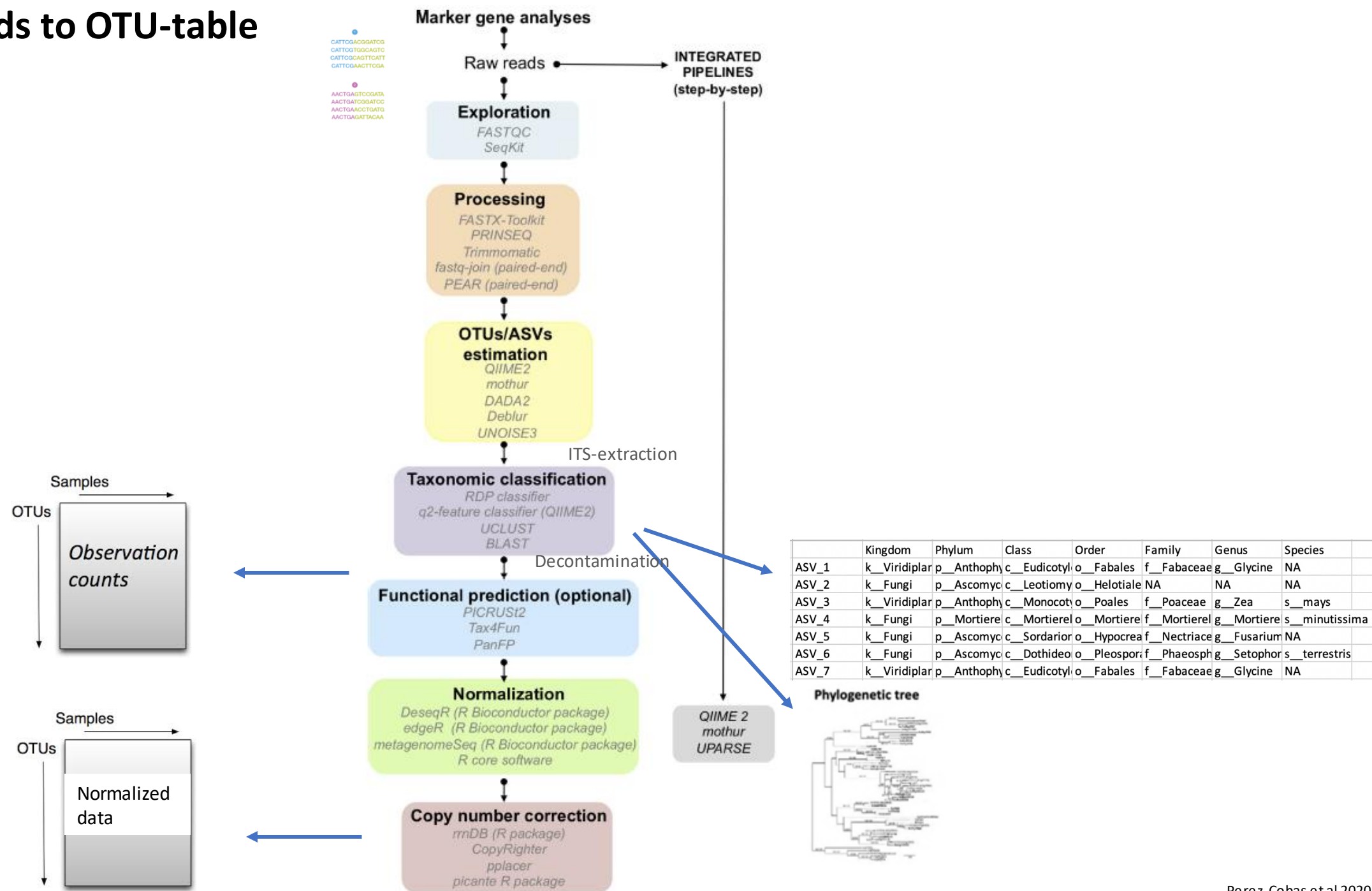
1. From raw-reads to OTU-table



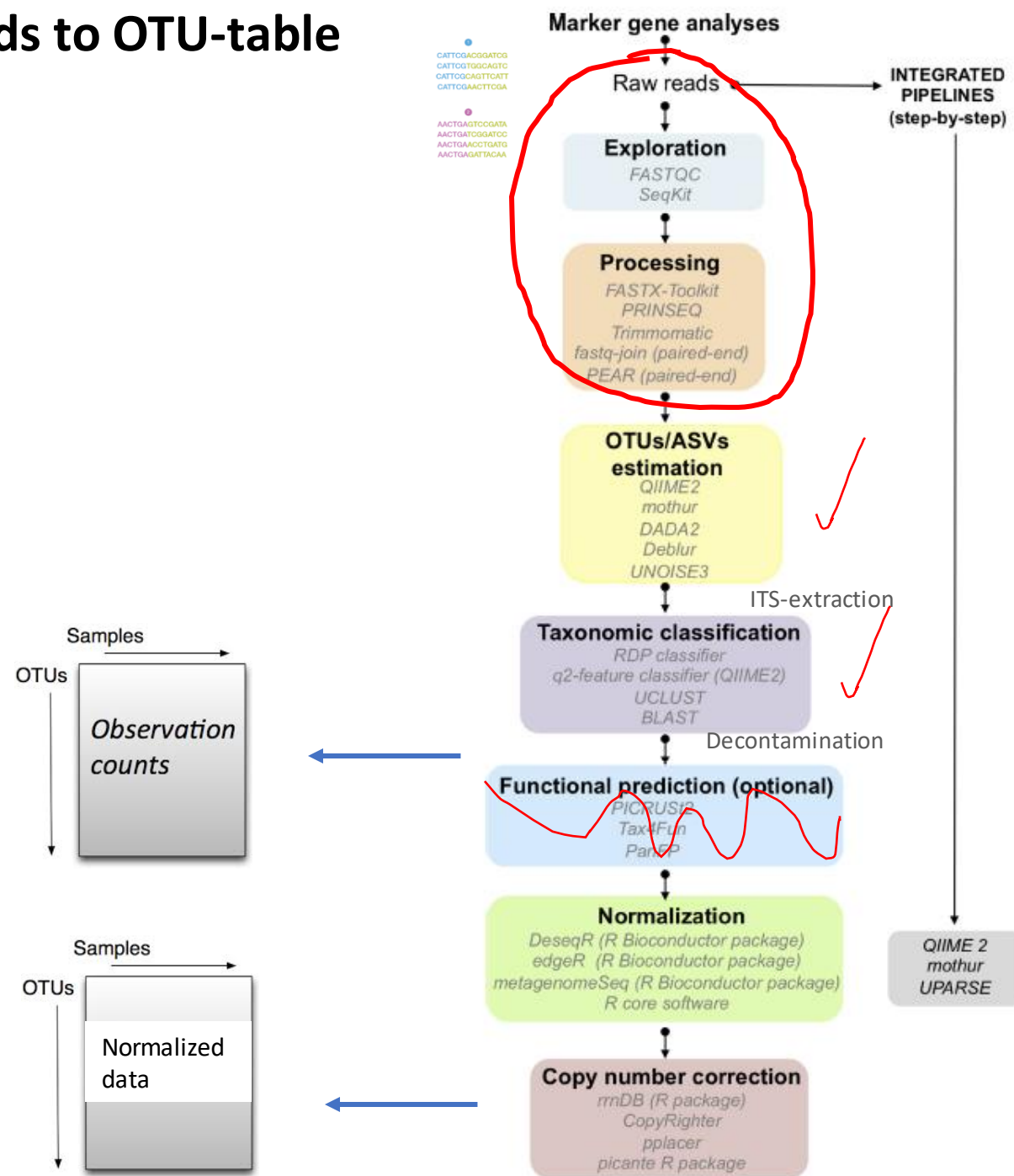
2. Statistical analysis (examples)



From raw-reads to OTU-table



From raw-reads to OTU-table



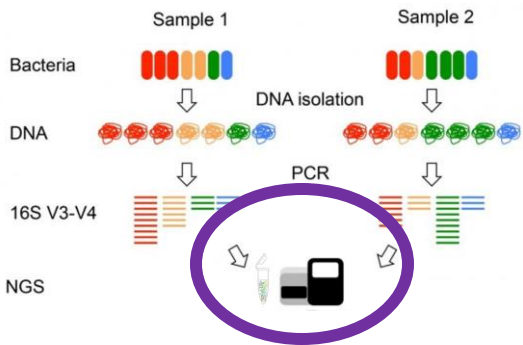
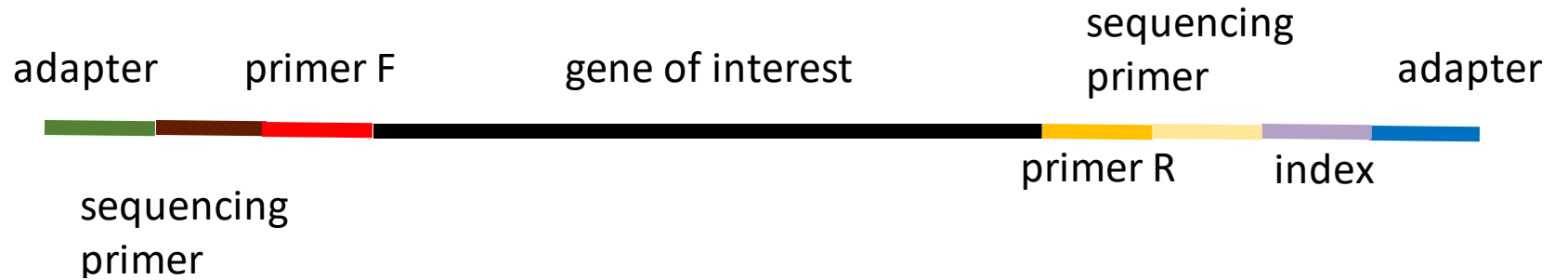
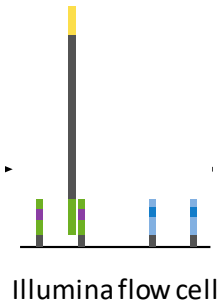
CATTGACCGATCG
 CATTGCTGGCAGTC
 CATTGCACTTCAT
 CATTGCAACTTGA

AACTGAGTCCGATA
 AACTGATCCGATCG
 AACTGAACTTGATG
 AACTGAGATTACAA

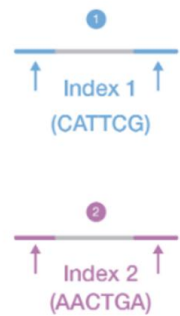
	Kingdom	Phylum	Class	Order	Family	Genus	Species	
ASV_1	k_Viridiplar	p_Anthophy	c_Eudicotyl	o_Fabales	f_Fabaceae	g_Glycine	NA	
ASV_2	k_Fungi	p_Ascomyc	c_Leotiomy	o_Helotiale	NA	NA	NA	
ASV_3	k_Viridiplar	p_Anthophy	c_Monocot	o_Poales	f_Poaceae	g_Zea	s_mays	
ASV_4	k_Fungi	p_Mortiere	c_Mortiere	o_Mortiere	f_Mortiere	g_Mortiere	s_minutissima	
ASV_5	k_Fungi	p_Ascomyc	c_Sordarior	o_Hypocrea	f_Nectriace	g_Fusarium	NA	
ASV_6	k_Fungi	p_Ascomyc	c_Dothideo	o_Pleospor	f_Phaeosph	g_Setophor	s_terrestris	
ASV_7	k_Viridiplar	p_Anthophy	c_Eudicotyl	o_Fabales	f_Fabaceae	g_Glycine	NA	



Structure of the amplicon (sequencing read)



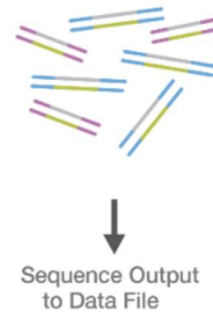
A
Library Preparation



B
Pool



C
Sequence

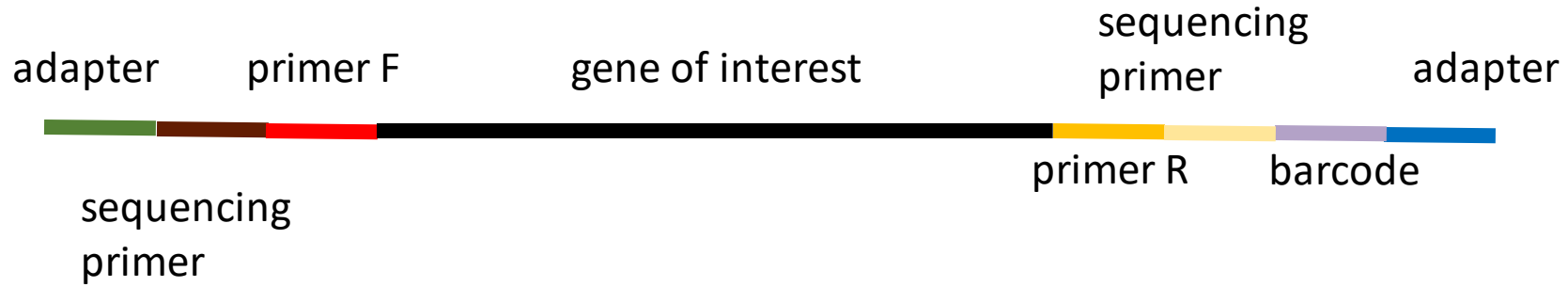
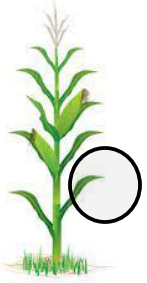


CATTGACGGATCG
AACTGAGTCCGATA
AACTGATCGGATCC
CATTCGTGGCAGTC
AACTGAACCTGATG
AACTGAGATTACAA
CATTGCGAGTTCATT
CATTGAACTTCGA

D
Demultiplex



Exploration and processing of reads



R1

R1

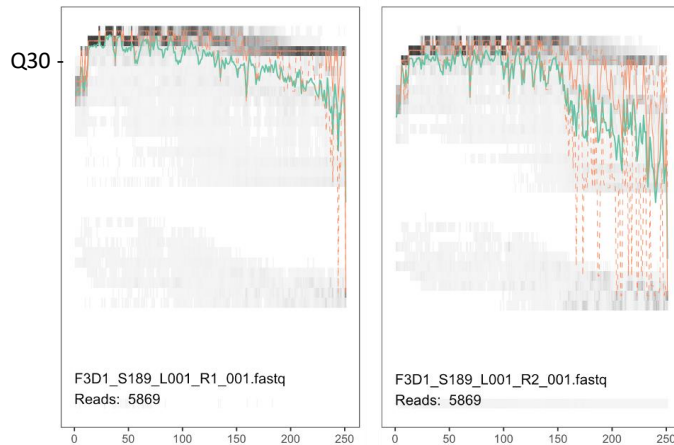
R2

fastq file

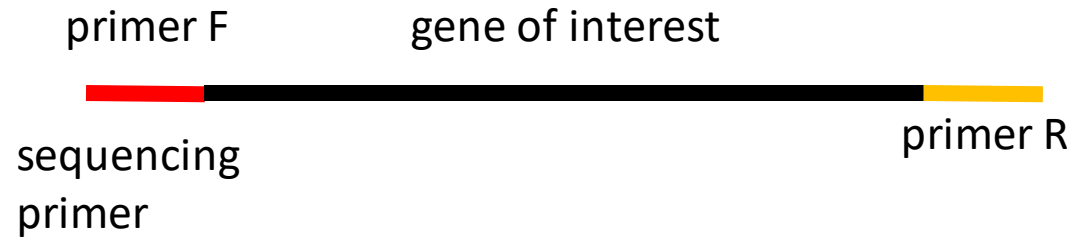
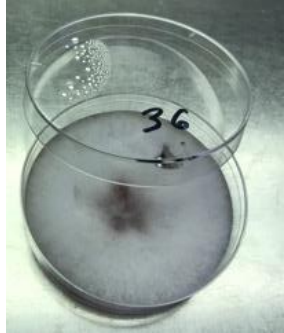
```

Identifier ● @SRR566546.970 HWUSI-EAS1673_11067_FC7070M:4:1:2299:1109 length=50
Sequence ● TTGCGTGCCTATCATTITAGTGCCTGTGAGGTGGAGATGTGAGGATCAGT
'+ sign ● +
Quality scores ● hhhhhhhhhghghghhhhhfhhhhfffffe'ee['X]b[d[ed['Y['Y
Identifier ● @SRR566546.971 HWUSI-EAS1673_11067_FC7070M:4:1:2374:1108 length=50
Sequence ● GATTTGTATGAAAGTATACAACATAAACTGCAAGTGGATCAGAGTAAGTC
'+ sign ● +
Quality scores ● hhhhgfhcghghggfcffdhfehhhcchhdhahehffffde'bVd
    
```

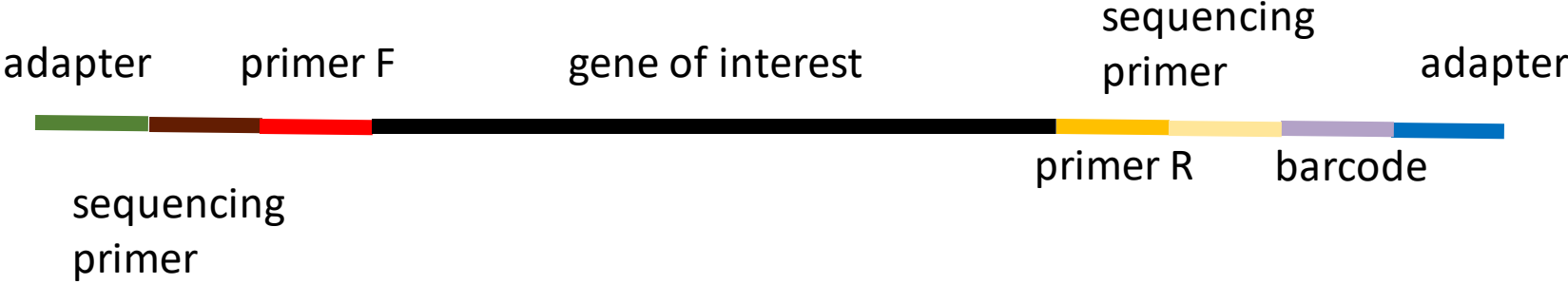
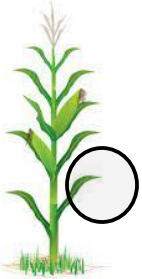
FASTQ – quality profile plot



Comparison with Sanger sequencing



Exploration and processing of reads

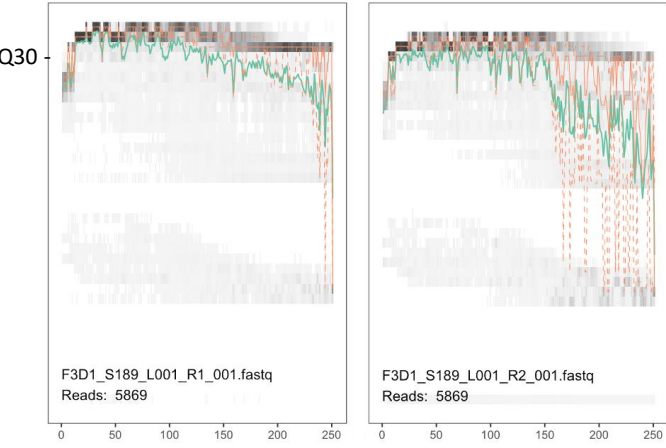


fastq file

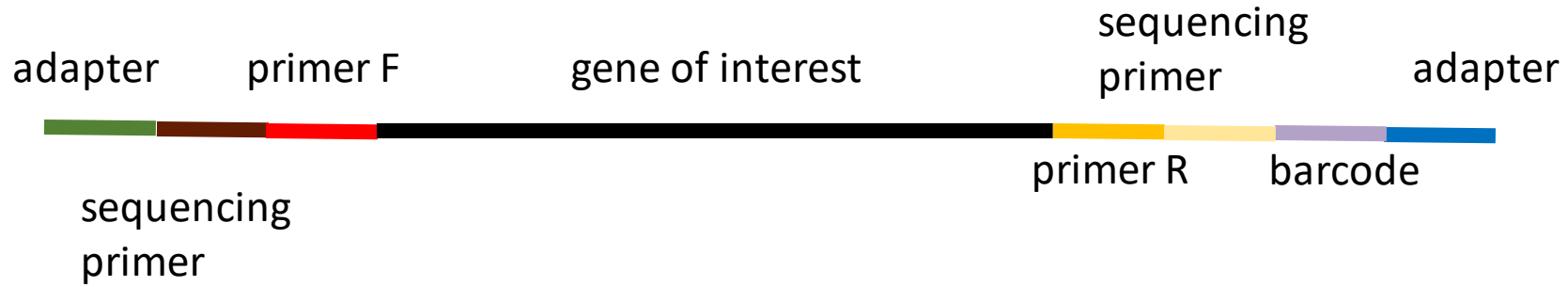
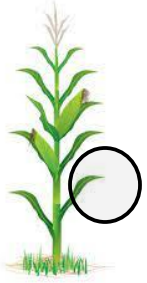
```

Identifier ● @SRR566546.970 HWUSI-EAS1673_11067_FC7070M:4:1:2299:1109 length=50
Sequence ● TTGCGTGCCTATCATTITAGTGCCTGTGAGGTGGAGATGTGAGGATCAGT
'+ sign ● +
Quality scores ● hhhhhhhhhghghghhhhhfhhhhfffffe'ee['X]b[d[ed['Y['Y
Identifier ● @SRR566546.971 HWUSI-EAS1673_11067_FC7070M:4:1:2374:1108 length=50
Sequence ● GATTTGTATGAAAGTATACAACATAAACTGCAAGGTGGATCAGAGTAAGTC
'+ sign ● +
Quality scores ● hhhhgfhcghghggfcffdhfehhhcchhdhahehffffde'bVd
    
```

FASTQ – quality profile plot



Exploration and processing of reads



R1

R1

R2



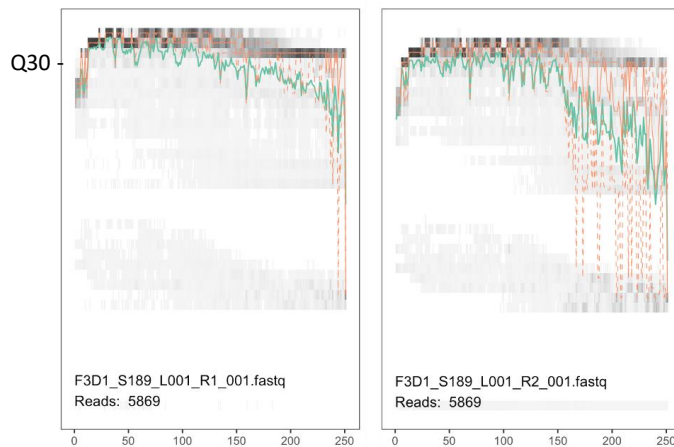
Taxonomic unit of analysis

fastq file

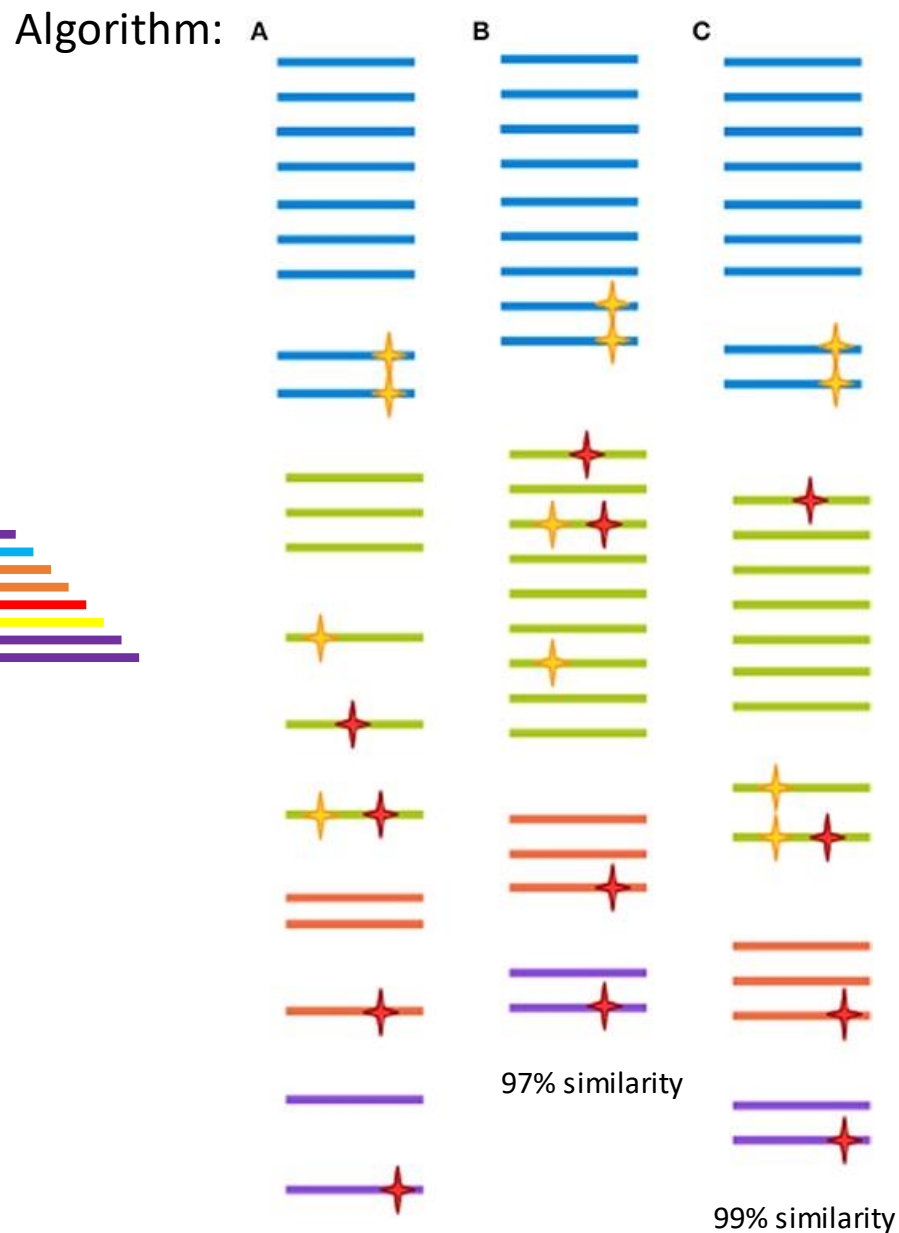
```

Identifier ● @SRR566546.970 HWUSI-EAS1673_11067_FC7070M:4:1:2299:1109 length=50
Sequence ● TTGCTGCCTATCATTITAGTGCCTGTGAGGTGGAGATGTGAGGATCAGT
'+ sign ● +
Quality scores ● hhhhhhhhhghghghhhhhfhhhhfffffe'ee['X]b[d[ed['Y['Y
Identifier ● @SRR566546.971 HWUSI-EAS1673_11067_FC7070M:4:1:2374:1108 length=50
Sequence ● GATTTGTATGAAAGTATACAACATAAACTGACAGGTGGATCAGAGTAAGTC
'+ sign ● +
Quality scores ● hhhhgfhcghghggfcffdhfehhhcchhdhahehffffde'bVd
    
```

FASTQ – quality profile plot



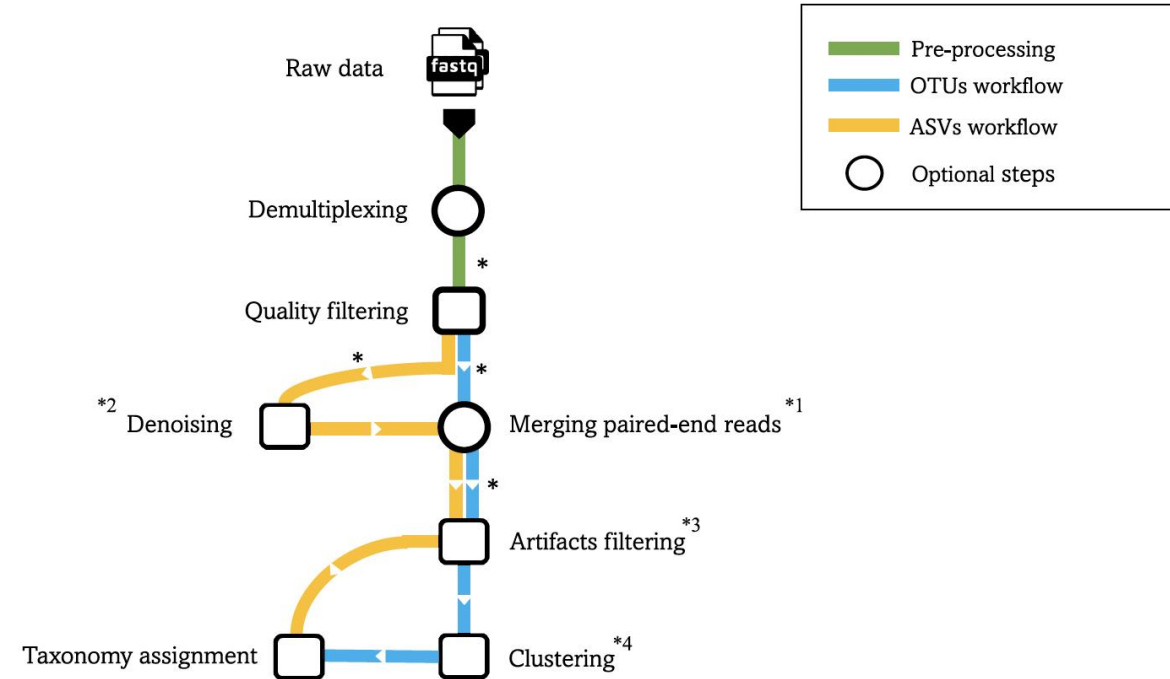
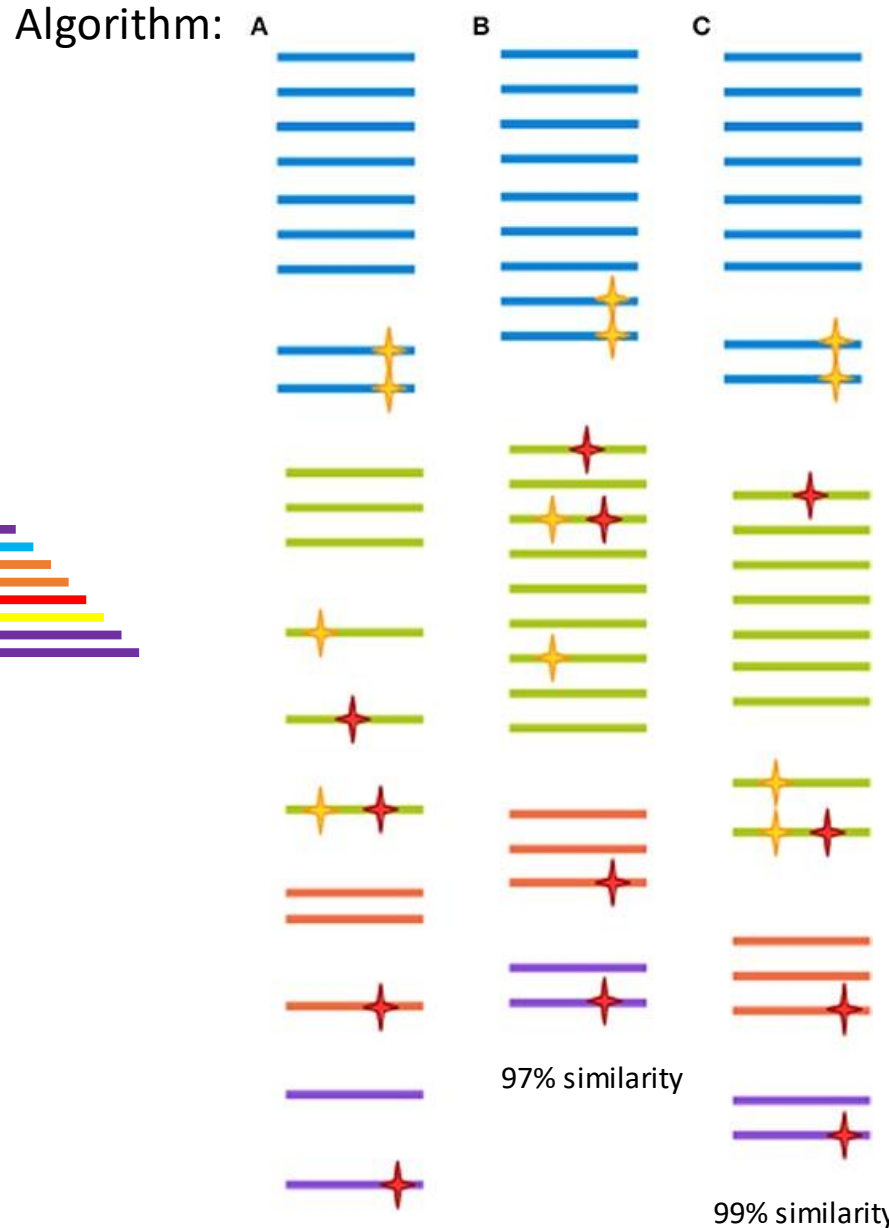
Estimation of the taxonomic unit of analysis (*De novo*)



yellow star: single nucleotide difference
red star: sequencing error

Single nucleotide difference

Estimation of the taxonomic unit of analysis (*De novo*)



yellow star: single nucleotide difference
red star: sequencing error

Single nucleotide difference

99% similarity

Hakimzadeh et al 2023

Hugerth and Andersson 2017

Denoising and error correction

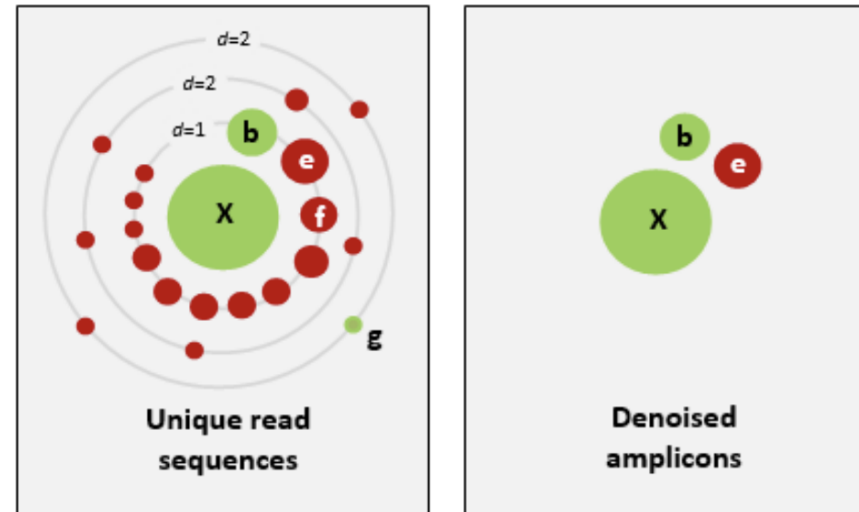


Figure 1. Schematic of the UNOISE2 denoising strategy. The left panel shows the neighborhood close to a high-abundance unique read sequence **X**, grouped by the number of sequence differences (d). Dots are unique sequences, the size of a dot indicates its abundance. Green dots are correct biological sequences; red dots have one or more errors. Neighbors with small numbers of differences and small abundance compared to **X** are predicted to be bad reads of **X**. The right panel shows the denoised amplicons. Here, **X** and **b** were correctly predicted, **e** is an error with anomalously high abundance that was wrongly predicted to be correct, **f** is an error that was correctly discarded but has an abundance almost high enough to be a false positive, and **g** is a low-abundance correct amplicon that was wrongly discarded. The abundances of **b**, **e**, and **f** are similar, illustrating the fundamental challenge in denoising: how to set an abundance threshold that distinguishes correct sequences from errors. *zOTUs (zero radius) or ESV (exact sequence variant)

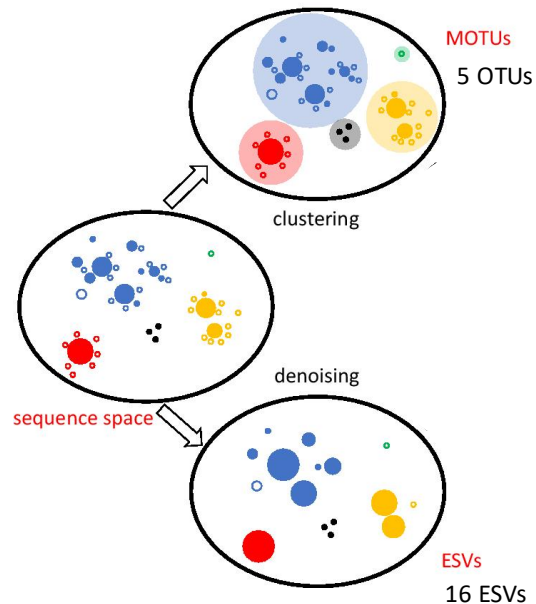
ASVs or OTUs?

~~ASVs or OTUs?~~

Does the pipeline incorporate denoising and error correction?

“ASVs are identical denoised reads with as few as 1 base pair difference between variants, representing an inference of the biological sequences prior to amplification and sequencing errors (Callahan et al., 2017)”.

OTUs represent clusters of sequences based on a specified similarity threshold. One sequence representative is chosen for further analysis.

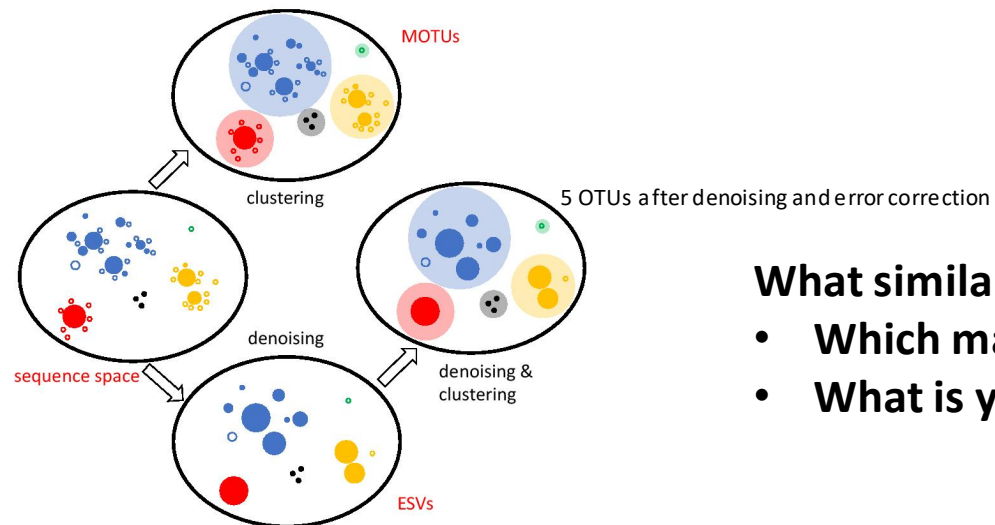


ASVs or OTUs?

Does the pipeline incorporate denoising and error correction?

“ASVs are identical denoised reads with as few as 1 base pair difference between variants, representing an inference of the biological sequences prior to amplification and sequencing errors (Callahan et al., 2017)”.

OTUs represent clusters of sequences based on a specified similarity threshold. One sequence representative is chosen for further analysis.



What similarity threshold to cluster with?

- **Which marker gene did you use?**
- **What is your desired taxonomic resolution?**

Relationship between gene region variability, % similarity, and 'genome splitting'

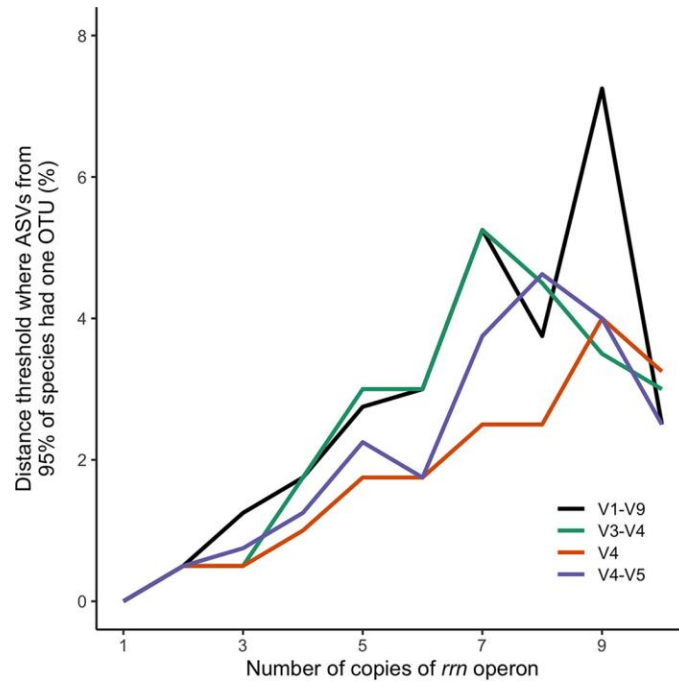


FIG 1 The distance threshold required to prevent the splitting of genomes into multiple OTUs increased as the number of *rrm* operons in the genome increased. Each line represents the median distance threshold for each region of the 16S rRNA gene that is required for 95% of the genomes with the indicated number of *rrm* operons to cluster their ASVs to a single OTU. The median distance threshold was calculated across 100 randomizations in which one genome was sampled from each species. Only those numbers of *rrm* operons that were found in more than 100 species are included.

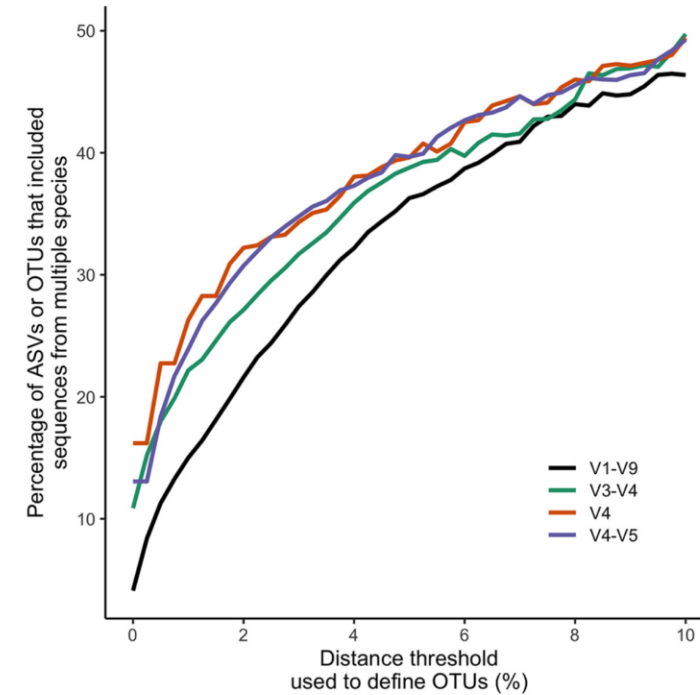
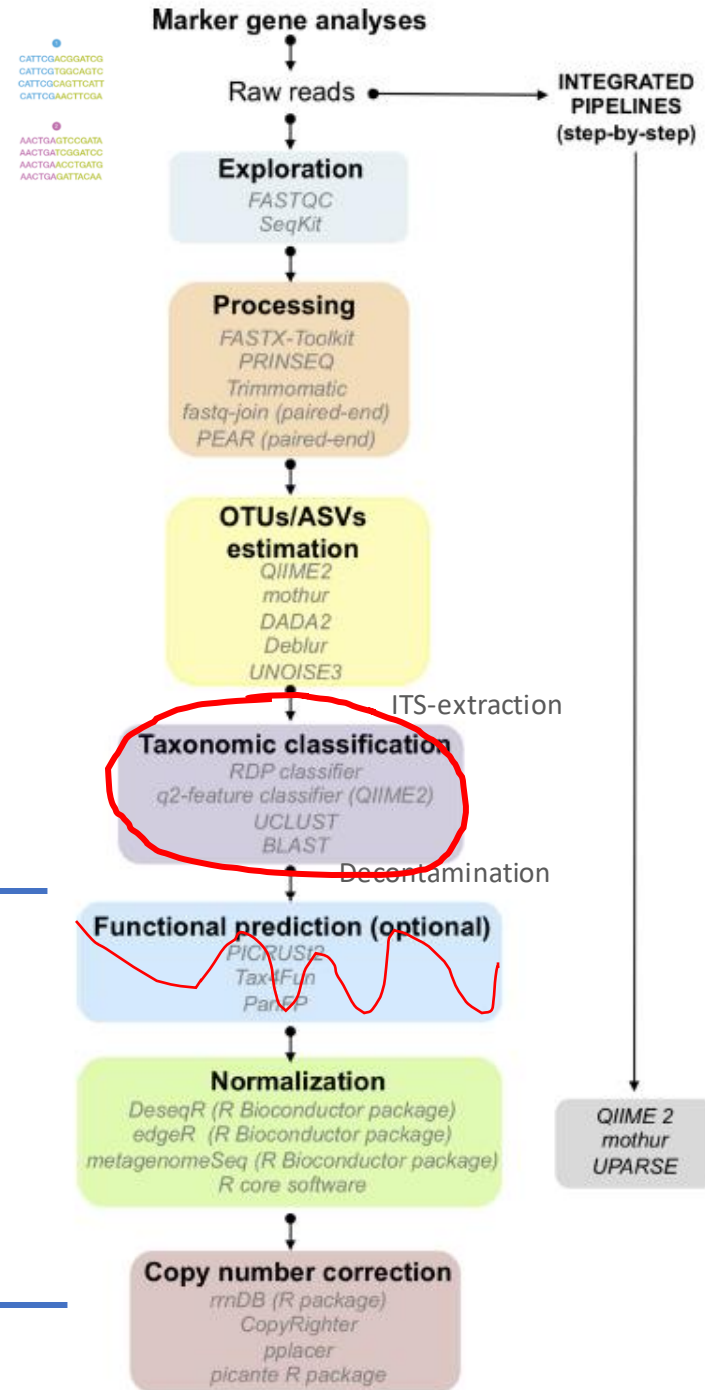
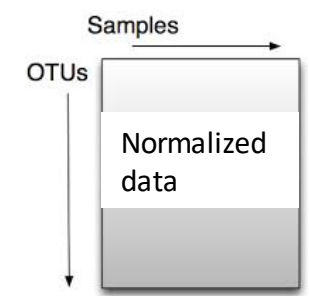
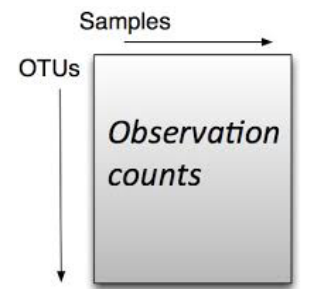


FIG 2 As the distance threshold used to define an OTU increased, the percentages of ASVs and OTUs representing multiple species increased. These data represent the median fractions for both measurements across 100 randomizations. In each randomization, one genome was sampled from each species.

From raw-reads to OTU-table



● CATTGACCGATCG
 CATTGCTGGCAGTC
 CATTGACGTTGAT
 CATTGGAACCTGGA
 ● AACTGAGTCCGATA
 AACTGATCCGATCG
 AACTGAACTCGATG
 AACTGAGATTACAA



	Kingdom	Phylum	Class	Order	Family	Genus	Species	
ASV_1	k_Viridiplar	p_Anthophyc	c_Eudicotyl	o_Fabales	f_Fabaceae	g_Glycine	NA	
ASV_2	k_Fungi	p_Ascomyc	c_Leotiomy	o_Helotiale	NA	NA	NA	
ASV_3	k_Viridiplar	p_Anthophyc	c_Monocoty	o_Poales	f_Poaceae	g_Zea	s_mays	
ASV_4	k_Fungi	p_Mortiere	c_Mortierel	o_Mortiere	f_Mortierel	g_Mortiere	s_minutissima	
ASV_5	k_Fungi	p_Ascomyc	c_Sordarior	o_Hypocrea	f_Nectriace	g_Fusarium	NA	
ASV_6	k_Fungi	p_Ascomyc	c_Dothideo	o_Pleospor	f_Phaeosph	g_Setophor	s_terrestris	
ASV_7	k_Viridiplar	p_Anthophyc	c_Eudicotyl	o_Fabales	f_Fabaceae	g_Glycine	NA	



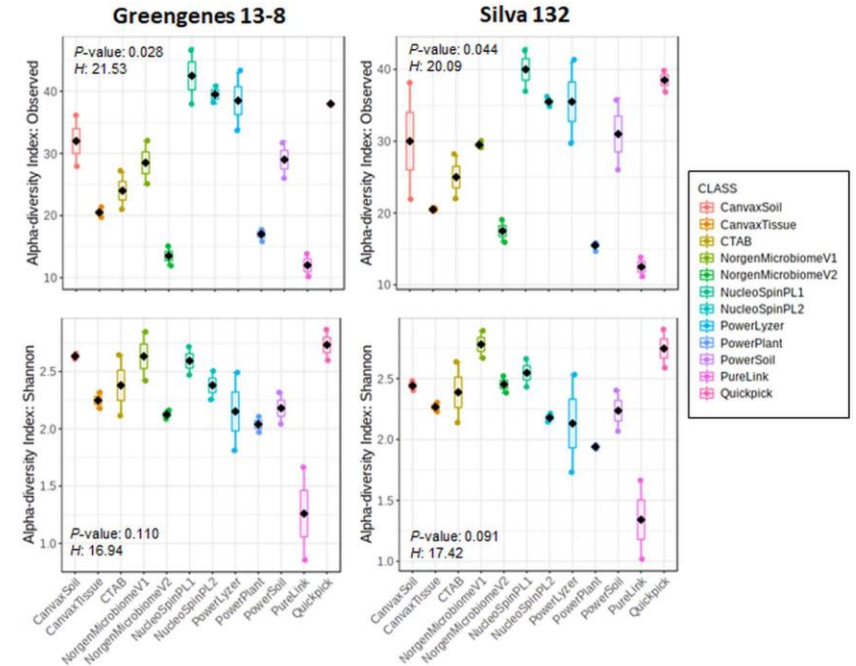
Taxonomic classification and databases

Example

```
>ASV1
GAGTTTGATCCTGGCTCAGGATGAACGCTGGCGGCGTGTCTTAACACATGCAAGTCGAACGGTGAAGCAG
GAGCTTGCTCTTGTGGATCAGTGGCGAACGGGTGAGTAACACGTGAGCAACCTGCCCCGAACCTGGGA
TAAGCGCTGGAAACGGCGTCTAATACTGGATATGCACCAGGGAGGCATCTTCACTGGTGGGAAAGATTTT
TTGGTTCGGGATGGGCTCGCGGCCTATCAGCTTGTGGTGAGGTAACGGCTCACCAAGGCGTCGACGGG
TAGCCGGCCTGAGAGGGTGACCGGCCACACTGGGACTGAGACACGGCCCAGACTCTACGGGAGGCAG
CAGTGGGGAATATTGCACAATGGGCGGAAGCCTGATGCAGCAACGCCG
```

Classification

- Bacteria
- Firmicutes
- Bacilli
- Bacillales
- Bacillaceae
- Bacillus
- Bacillus subtilis



*depends on the taxonomic resolution of your marker gene (and its length)

Taxonomic classification and databases

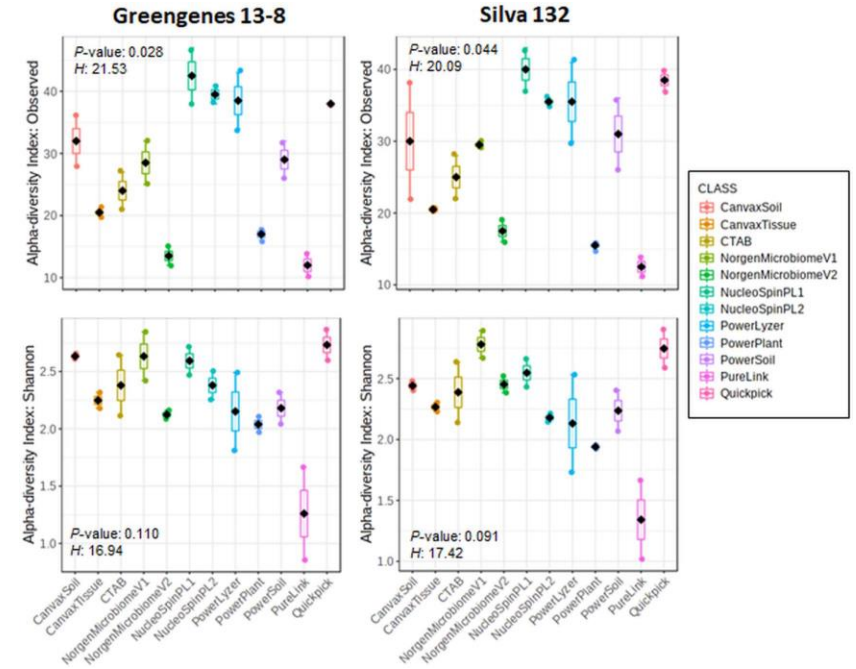
Example

```
>ASV1
GAGTTTGATCCTGGCTCAGGATGAACGCTGGCGGCGTGTCTTAACACATGCAAGTCGAACGGTGAAGCAG
GAGCTTGCTCTTGTGGATCAGTGGCGAACGGGTGAGTAACACGTGAGCAACCTGCCCCGAACCTGGGA
TAAGCGCTGGAAACGGCGTCTAATACTGGATATGCACCAGGGAGGCATCTTCACTGGTGGGAAAGATTTT
TTGGTTCGGGATGGGCTCGCGGCCTATCAGCTTGTGGTGAGGTAACGGCTCACCAAGGCGTCGACGGG
TAGCCGGCCTGAGAGGGTGACCGGCCACACTGGGACTGAGACACGGCCCAGACTCTACGGGAGGCAG
CAGTGGGGAATATTGCACAATGGGCGGAAGCCTGATGCAGCAACGCCG
```

Classification

- Bacteria
- Firmicutes
- Bacilli
- Bacillales
- Bacillaceae
- Bacillus
- Bacillus subtilis

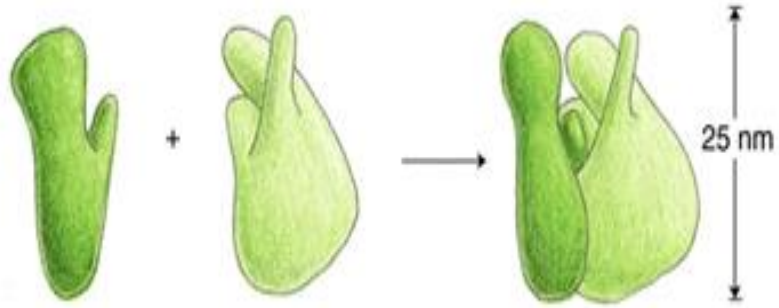
Tree building (all ASVs)



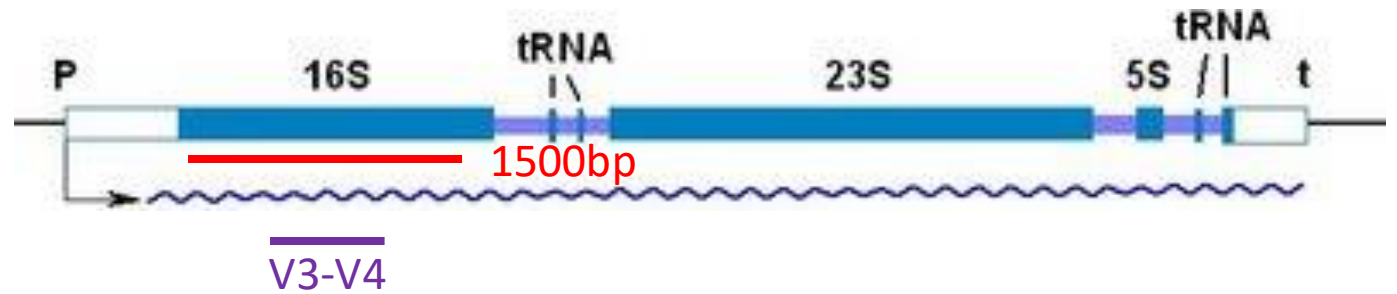
*depends on the taxonomic resolution of your marker gene (and its length)

Ribosomal markers as taxonomic barcodes

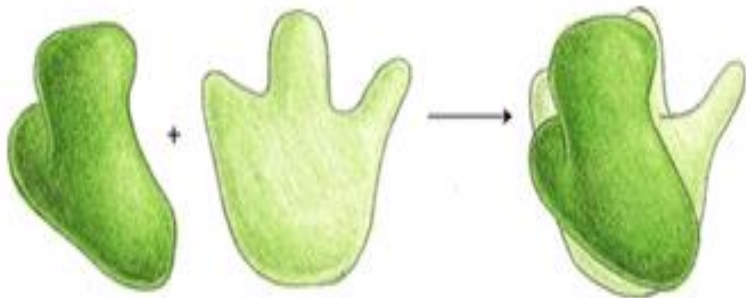
30S 50S



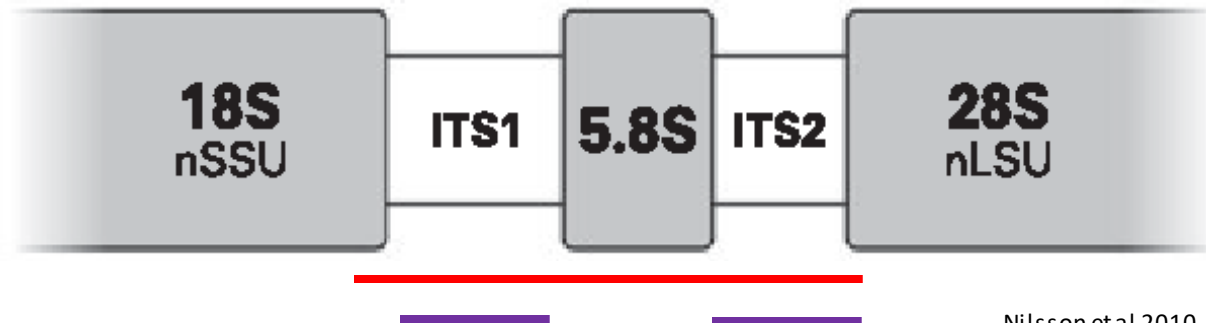
Bacteria (chloroplast, mitochondria)



40S 60S



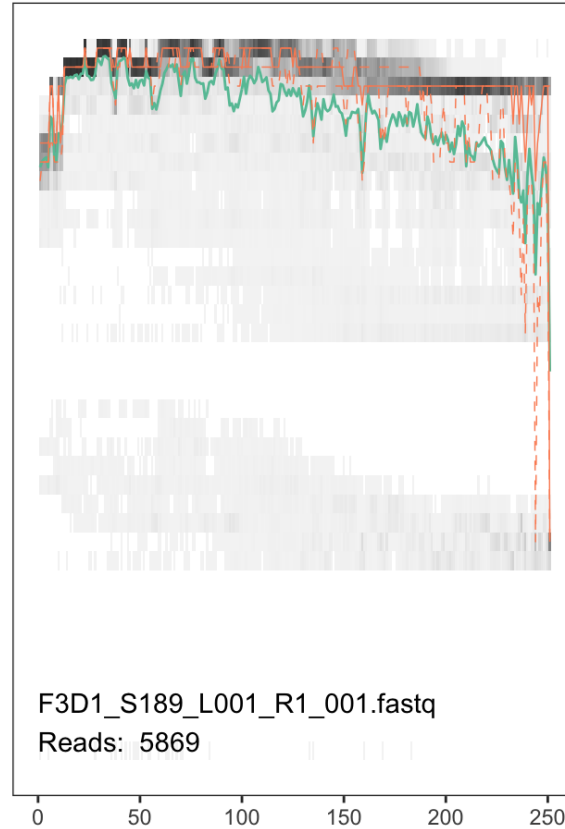
Eukarya (Fungi, Nematodes, Protists, Plants)



Things to look for in a pipeline

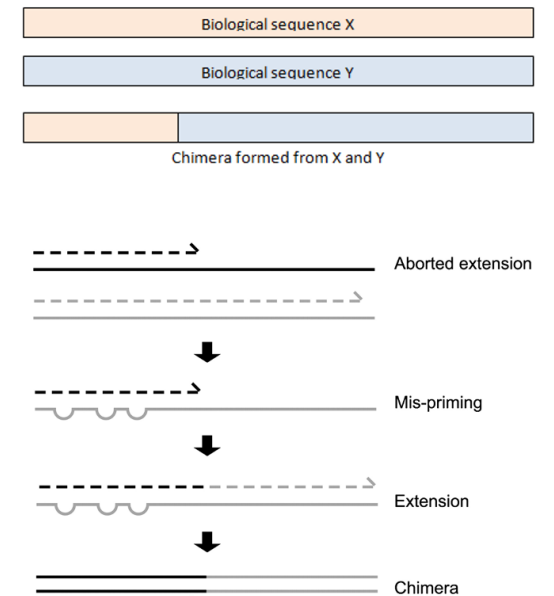
- Different quality control steps:
 - Sequence length (trimming)
 - Low quality reads (filtering)
 - Deal with sequencing errors (denoising/error correction)
 - Homopolymers
 - Chimeras
- OTU vs SV
- Deal with controls
- Documentation
- Format/compatibility with downstream analysis

FASTQ – quality profile plot



https://benjjneb.github.io/dada2/tutorial_1_6.html

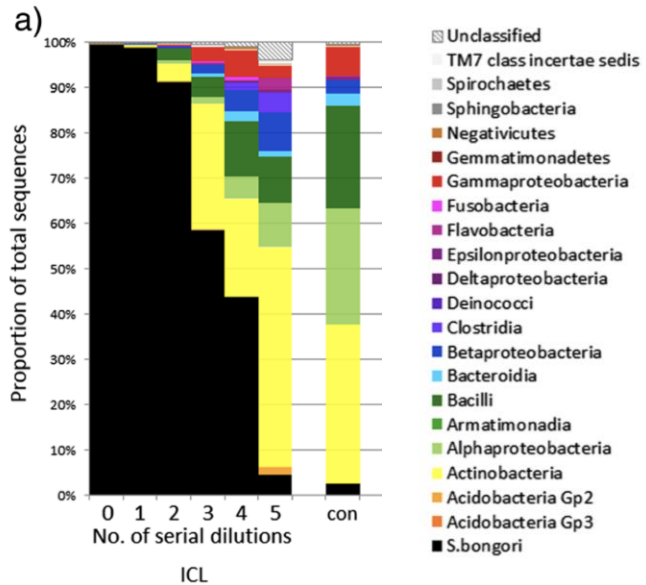
PCR and Chimeras



Bahram et al (2019) doi:[10.1111/1758-2229.12684](https://doi.org/10.1111/1758-2229.12684); Haas et al 2011

Negative controls and identifying contaminants

Contaminants in low biomass samples



<http://www.biomedcentral.com/1741-7007/12/87>

Davis et al. *Microbiome* (2018) 6:226
<https://doi.org/10.1186/s40168-018-0605-2>

Microbiome

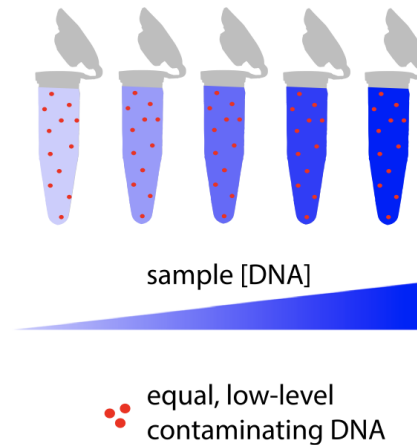
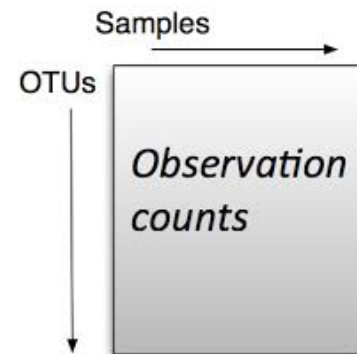
METHODOLOGY

Open Access

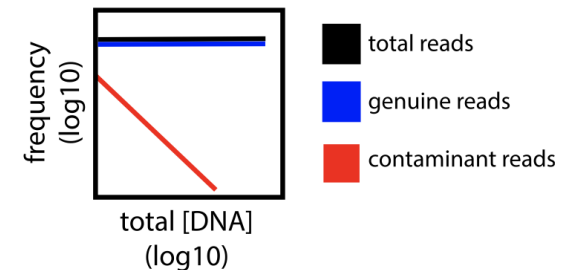
Simple statistical identification and removal of contaminant sequences in marker-gene and metagenomics data



Nicole M. Davis¹, Diana M. Proctor^{2,3}, Susan P. Holmes⁴, David A. Relman^{1,2,5} and Benjamin J. Callahan^{6,7*}



sequence equimolar amounts well-mixed total DNA



contaminant DNA correlates inversely with total DNA

Technical replicates or other standards

Phytobiomes • 2018 • 2:165-170

<https://doi.org/10.1094/PBIOMES-09-17-0041-R>



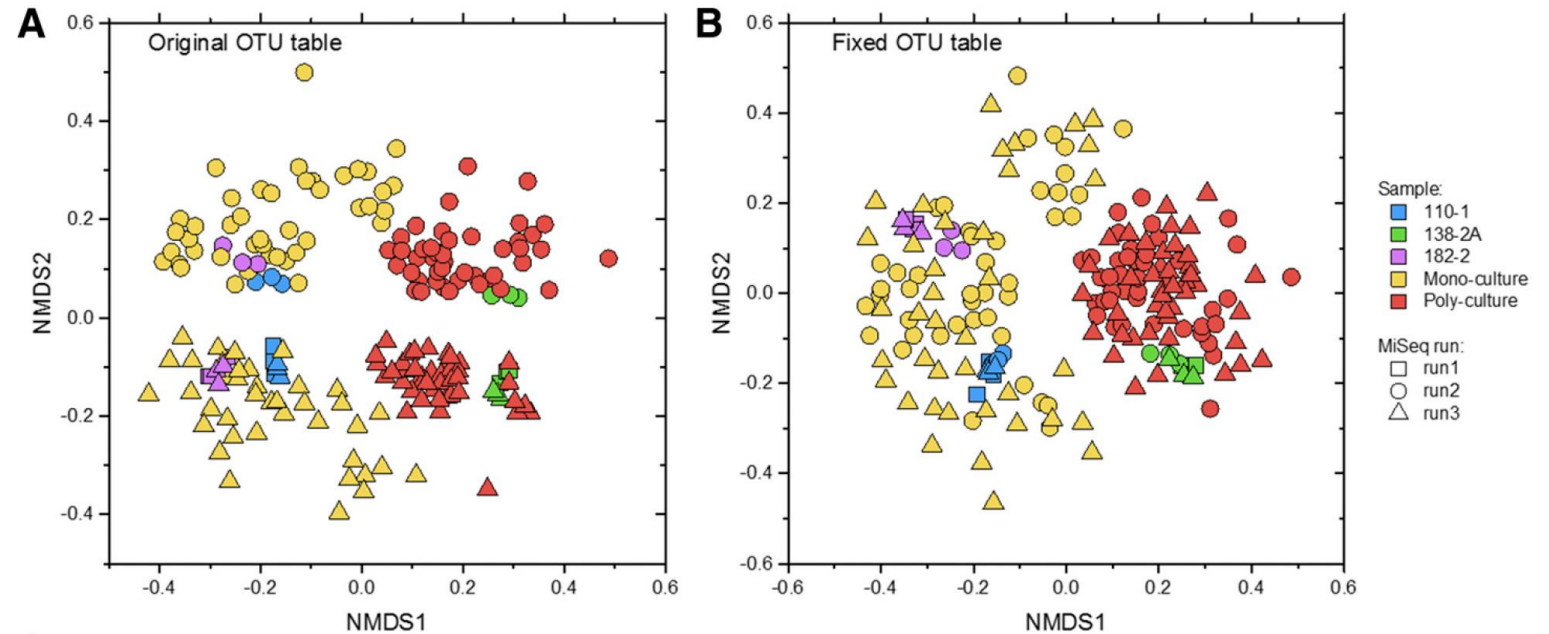
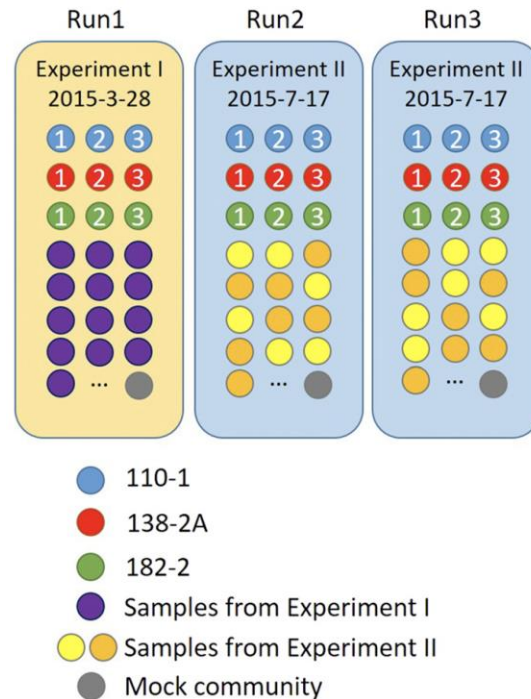
RESEARCH

e-Xtra*

Run-to-Run Sequencing Variation Can Introduce Taxon-Specific Bias in the Evaluation of Fungal Microbiomes

Zewei Song[†] and Dan Schlatter, Department of Plant Pathology, University of Minnesota, Saint Paul; Daryl M. Gohl, University of Minnesota Genomics Center, Minneapolis; and Linda L. Kinkel, Department of Plant Pathology, University of Minnesota, Saint Paul

Experiment design

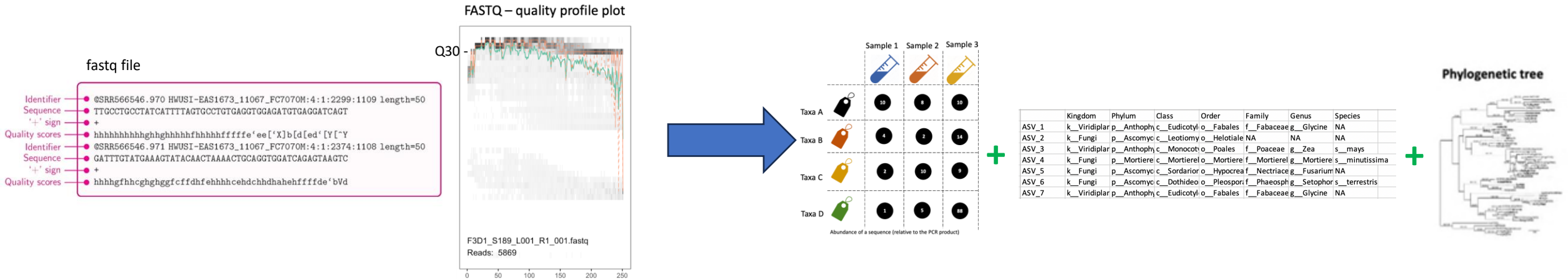


Mock communities and positive controls

Product	Catalog #	Composition	Format
Microbial Community Standard	D6300	Even Distribution	Microbial
Microbial Community DNA Standard	D6305/6306	Even Distribution	Isolated DNA
Microbial Community Standard II	D6310	Log Distribution	Microbial
Microbial Community DNA Standard II	D6311	Log Distribution	Isolated DNA
Spike-in Control I (High Microbial Load)	D6320/D6320-10	Even Distribution	Microbial
Spike-in Control II (Low Microbial Load)	D6321/D6321-10	Log Distribution	Microbial
HMW DNA Standard	D6322	Even Distribution	Isolated DNA
Gut Microbiome Standard	D6331	Staggered Abundance	Microbial

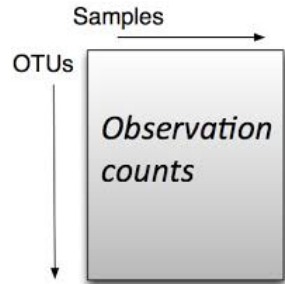
Metabarcoding processing and analysis pipelines

1. From raw-reads to OTU-table



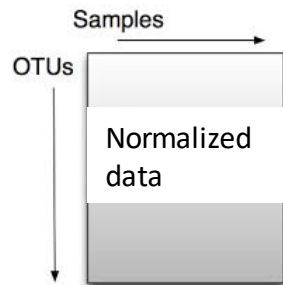
2. Statistical analysis

Analysing and interpreting your data

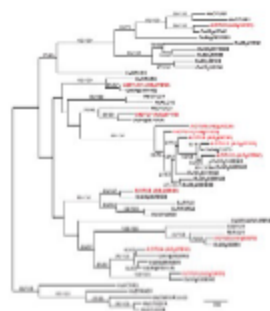


	A	B	C	D	E	F
1		sample01	sample01	sample02	sample02	sample03
2	OTU0001	528	68	1755	773	2167
3	OTU0002	138	68	559	2588	1198
4	OTU0003	36	533	673	351	815
5	OTU0004	1	2618	5	17	19
6	OTU0005	224	237	81	271	313

	Kingdom	Phylum	Class	Order	Family	Genus	Species
ASV_1	k_Viridiplar	p_Anthophy	c_Eudicotyl	o_Fabales	f_Fabaceae	g_Glycine	NA
ASV_2	k_Fungi	p_Ascomyc	c_Leotiomy	o_Helotiale	NA	NA	NA
ASV_3	k_Viridiplar	p_Anthophy	c_Monocot	o_Poales	f_Poaceae	g_Zea	s_mays
ASV_4	k_Fungi	p_Mortiere	c_Mortiere	o_Mortiere	f_Mortiere	g_Mortiere	s_minutissima
ASV_5	k_Fungi	p_Ascomyc	c_Sordarior	o_Hypocrea	f_Nectriace	g_Fusarium	NA
ASV_6	k_Fungi	p_Ascomyc	c_Dothideo	o_Pleospori	f_Phaeosph	g_Setophor	s_terrestris
ASV_7	k_Viridiplar	p_Anthophy	c_Eudicotyl	o_Fabales	f_Fabaceae	g_Glycine	NA



Phylogenetic tree



Metadata

Phytobiomes Journal • 2020 • 4:115-121

<https://doi.org/10.1094/PBIOMES-09-19-0051-P>

OPEN ACCESS

Phytobiomes Journal

phytobiomesjournal.org

APS

A Transdisciplinary Journal of Sustainable Plant Productivity

PERSPECTIVE

Community-Driven Metadata Standards for Agricultural Microbiome Research

J. P. Dundore-Arias,^{1,†} E. A. Eloe-Fadrosh,² L. M. Schriml,³ G. A. Beattie,⁴ F. P. Brennan,⁵ P. E. Busby,⁶ R. B. Calderon,⁷ S. C. Castle,⁸ J. B. Emerson,⁹ S. E. Everhart,¹⁰ K. Eversole,¹¹ K. E. Frost,¹² J. R. Herr,¹³ A. I. Huerta,¹⁴ A. S. Iyer-Pascuzzi,¹⁵ A. K. Kalii,¹⁶ J. E. Leach,¹⁷ J. Leonard,¹⁸ J. E. Maul,¹⁹ B. Prithiviraj,²⁰ M. Potrykus,²¹ N. R. Redekar,²² J. A. Rojas,²³ K. A. T. Silverstein,²⁴ D. J. Tomso,²⁵ S. G. Tringe,²⁶ B. A. Vinatzer,²⁷ and L. L. Kinkel²⁸

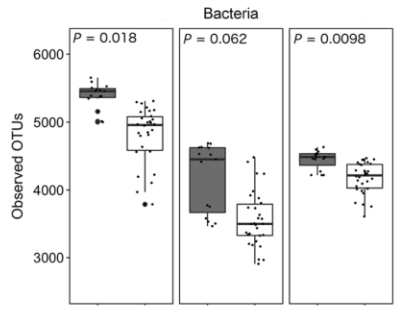
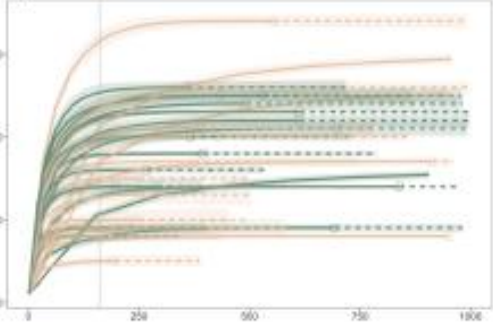


*phyloseq object

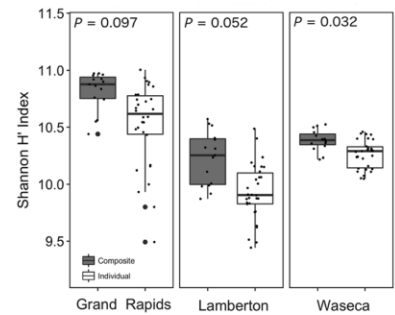
Which analyses to apply?

Alpha-diversity: within a community

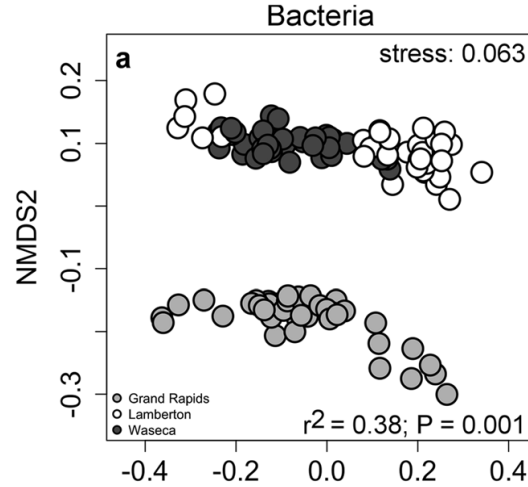
a) Root bacteria



Richness
Evenness
Phylogenetic diversity

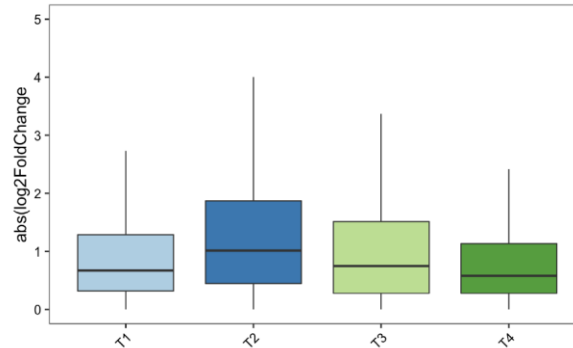


Beta-diversity: between communities

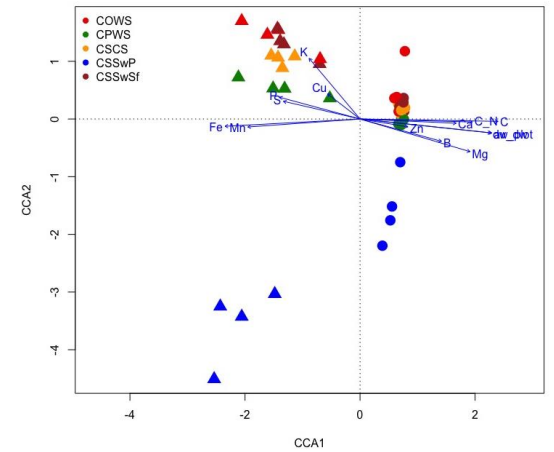
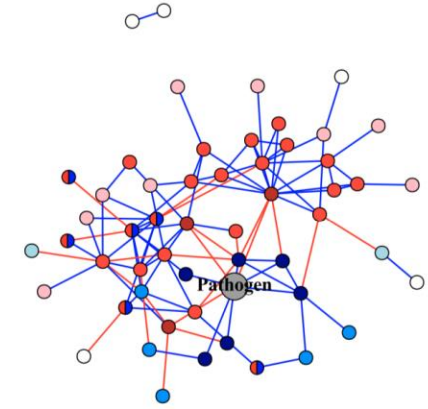


Similarity/dissimilarity indices
E.g: Bray-Curtis, Unifrac...

Differential abundance tests



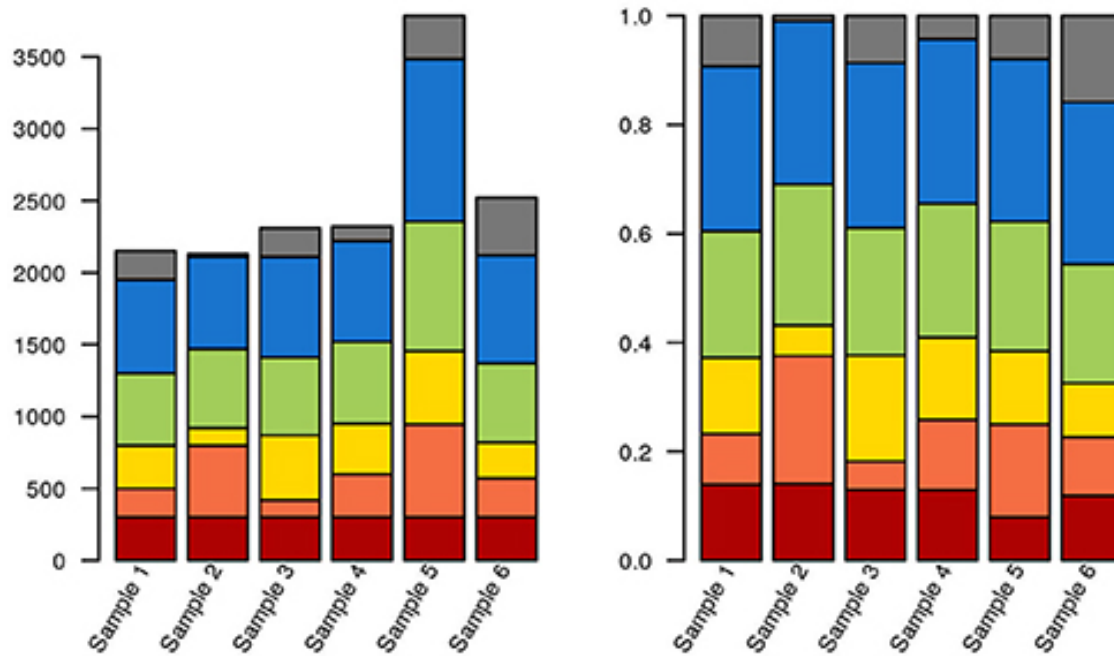
Co-occurrence and interactions between microbes; and correlation with environmental variables



To normalize or not?

What is your question and hypothesis?

McMurdie, Paul J., and Susan Holmes. "Waste Not, Want Not: Why Rarefying Microbiome Data Is Inadmissible." *PLoS Comput Biol* 10, no. 4 (2014): e1003531.



Weiss, Sophie, Zhenjiang Zech Xu, Shyamal Peddada, Amnon Amir, Kyle Bittinger, Antonio Gonzalez, Catherine Lozupone, et al. "Normalization and Microbial Differential Abundance Strategies Depend upon Data Characteristics." *Microbiome* 5, no. 1 (March 3, 2017): 27. <https://doi.org/10.1186/s40168-017-0237-y>.

Schloss PD. 2024. Waste not, want not: revisiting the analysis that called into question the practice of rarefaction. *mSphere* 9:e00355-23. <https://doi.org/10.1128/msphere.00355-23>

Variation in RNA operon copy number

rrndb: the Ribosomal RNA Operon Copy Number Database



Table 1.

Intra-genomic 16S rRNA variability for Bacteria and Archaea with full-genome sequence availability

Organism	No. rRNA ^a operons	Diff. (nt) ^b	% difference ^c
<i>Aquifex aeolicus</i> VF5	2	–	–
<i>Bacillus subtilis</i> ATCC 23857	10	1–15	0.97
<i>Campylobacter jejuni</i> ATCC 700819	3	–	–
<i>Deinococcus radiodurans</i> ATCC 13939	3	0–2	0.13
<i>Escherichia coli</i> ATCC 10798	7	0–19	1.23
<i>Haemophilus influenzae</i> ATCC 51907	6	–	–
<i>Helicobacter pylori</i> 26695	2	–	–
<i>Methanococcus jannaschii</i> DSMZ 2661	2	3	0.20
<i>Methanococcus thermoautotrophicum</i> ATCC 29096	2	2	0.14
<i>Neisseria meningitidis</i> MC 58	4	–	–
<i>Treponema pallidum</i> ATCC 25870	2	–	–
<i>Ureaplasma urealyticum</i> serovar 3	2	1	0.07
<i>Vibrio cholerae</i> ATCC 39315	8	0–14	0.91
<i>Xyella fastidiosa</i> 9a5c	2	–	–

^aNumber of rRNA operons per genome.

^bPairwise difference range between 16S rRNA genes per genome.

^cPairwise difference range between 16S rRNA genes per genome calculated as a percentage. –, no nucleotide differences.

