

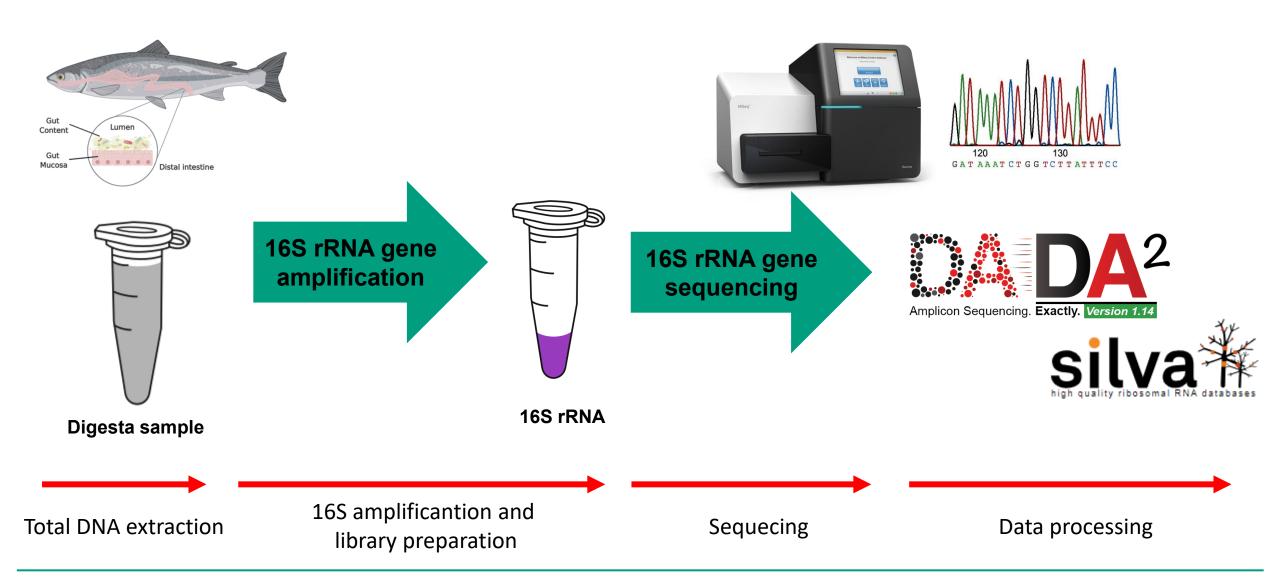
16S data processing and analysis

Dr. Sérgio Rocha

24.06.2025

General approach





Orion.nmbu.no

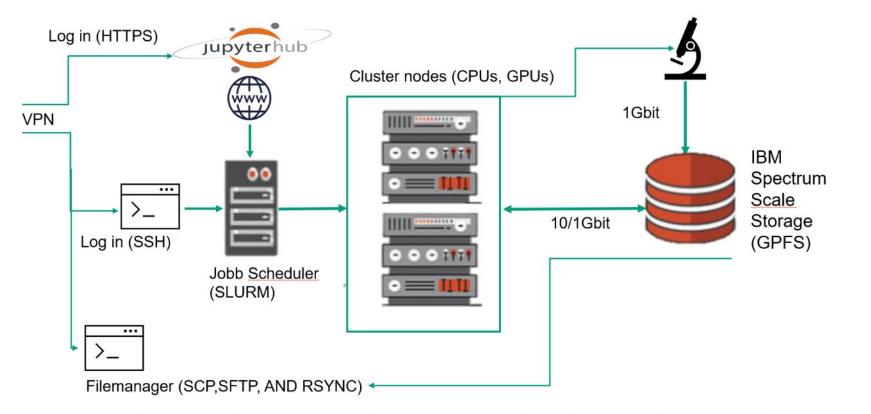


Orion HPC

Computing problems, such as genome sequence analysis, computational chemistry, simulation of universe activity, etc., require a different kind of computing power. These computing problems may need thousands, millions or billions of instructions to be completed

Main configuration of the Orion HPC

The Orion HPC system currently consists of 1680 processor cores with more than 12 terabytes of RAM and 1 petabyte of storage accessible on a 10/1 Gbit network via NFS. The Orion HPC's computation nodes utilize various processors with Centos Linux 7.9 as an operating environment.



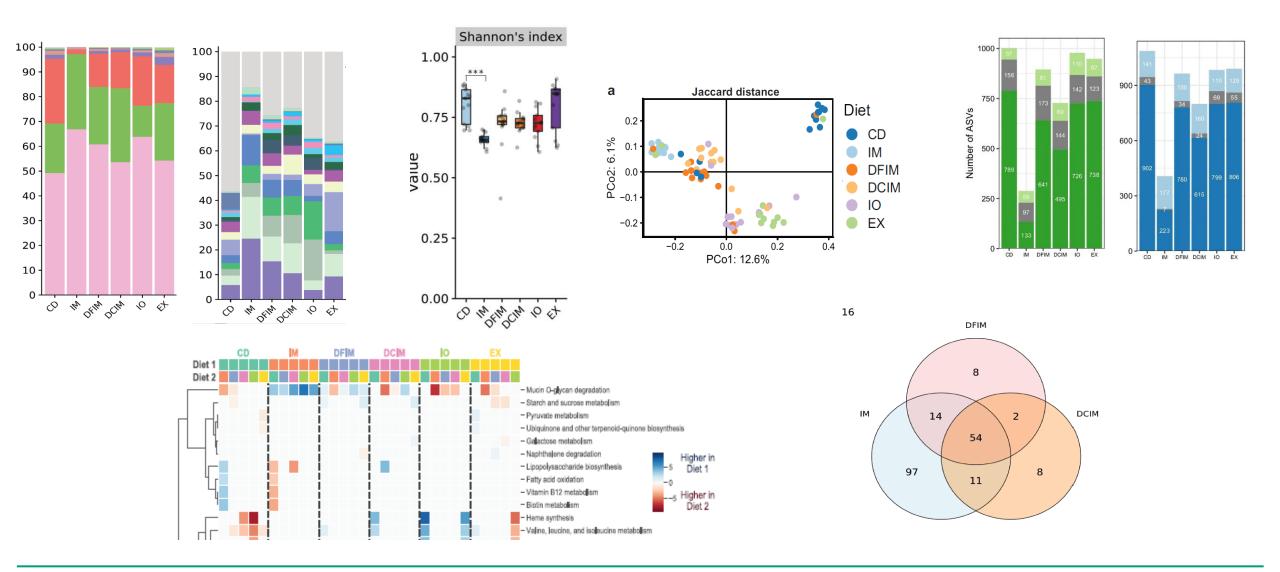


Orion HPC provides a supercomputing environment with thousands of cores to researchers and students to handle their computational problems. The following table summarizes the computing power capacity of Orion HPC nodes.

Examples







Some terminology



Amplicon: The resulting sequence of a targeted amplification of genetic material. Targeted meaning primers were used. (different from shotgut)

Marker gene: A gene that can be useful for delineating organisms, like a fingerprint (16S rRNA).

OUT: Operational Taxonomic Unit. OTUs are an artificial, arbitrary construct useful for grouping sequences together into units to help us summarize and analyze things, and they also intended to help deal with sequencing error intrinsic to the technology. \rightarrow 97% or 99% OTU clustering

ASV: <u>Amplicon Sequence Variant</u>. Resulting sequences from newer processing methodologies that attempt to take into account sequencing error rates and are believed to represent true biological sequences.

Moving towards using ASVs over OTUs is supported in recent publications.

Barcodes: sequences ligated to your individual samples' genetic material before they get all mixed together to be sequenced together. These barcodes are then unique to each sample, so you can afterwards identify which sequences came from which samples.

Demultiplex: step in processing amplicon sequence data where you would use the barcode information in order to know which sequences came from which samples after they had all been sequenced together.

OTUs or ASVs?



There are several approaches, and you will likely get very similar results regardless of which tool you use (as long as you make similar decisions when processing your sequences like minimum abundance filtering, ...).

Don't get too lost on trying to find the "best" tool for processing amplicon data.

There is a tendency on getting away from the traditional OTU approach and use more ASVs.

ASVs approach, is most often **more biologically meaningful** and a more useful unit beyond the current dataset than the traditional OTU clustering methods. "It allows to find new sequences".

If you are processing a new amplicon dataset, I would suggest to use DADA2, but if you have good support by using other pipeline, it is still a valid approach.

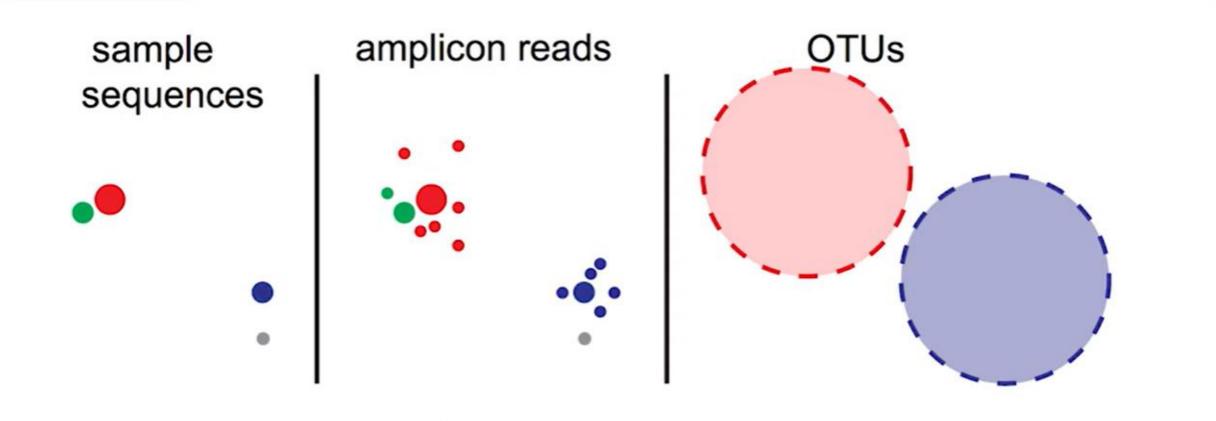
https://astrobiomike.github.io/amplicon/

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

Youtube: Bioinformatics Virtual Coordination Network

OTUs or ASVs?



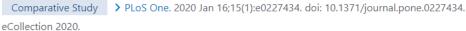


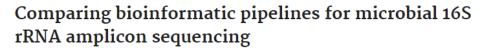
Bioinformatic pipelines











Andrei Prodan ¹, Valentina Tremaroli ², Harald Brolin ², Aeilko H Zwinderman ³, Max Nieuwdorp ¹, Evgeni Levin ¹ ⁴

Affiliations + expand

PMID: 31945086 PMCID: PMC6964864 DOI: 10.1371/journal.pone.0227434



but also usearch, vsearch, Minimum Entropy Decomposition, UNOISE ...

FULL TEXT LINKS

ACTIONS

OPEN ACCESS TO FULL TEXT

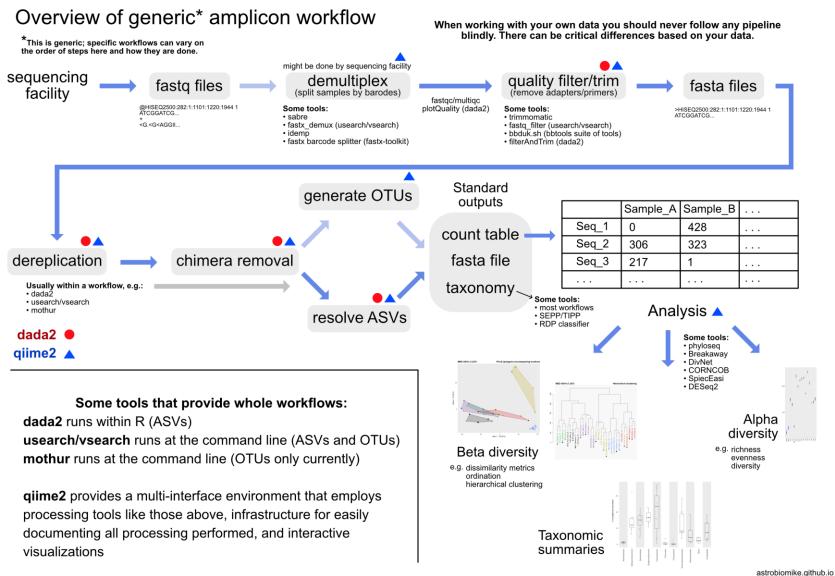
PLOS ONE

Cite

□ Collections

Amplicon workflow









nature methods

Explore content > About the journal > Publish with us > Subscribe

nature > nature methods > brief communications > article

Brief Communication | Published: 23 May 2016

DADA2: High-resolution sample inference from Illumina amplicon data

Benjamin J Callahan [™], Paul J McMurdie, Michael J Rosen, Andrew W Han, Amy Jo A Johnson & Susan P Holmes

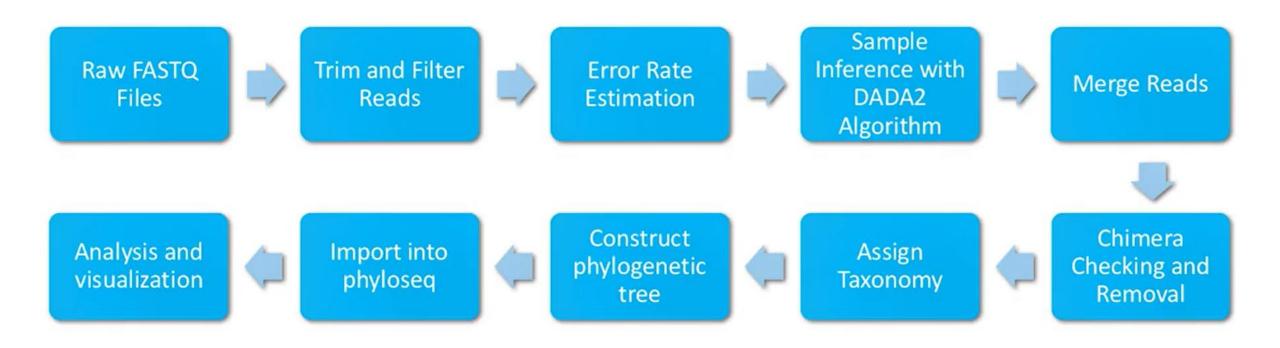
Nature Methods 13, 581–583 (2016) Cite this article

102k Accesses | 14k Citations | 114 Altmetric | Metrics

Bioinformatics



DADA2 Workflow



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

«FON pipeline»



- 1. Sequencing denoising
- 2. Pre-processing
- 3. Taxonomic Analysis
- 4. alpha-diversity
- 5. beta-diversity
- 6. Individual taxa comparisson
- 7. LEfSe
- 8. Metabolic pathways



Instal some basic packages

```
## Installation and loading of dada2

```{r}
if (!requireNamespace("BiocManager", quietly = TRUE))
 install.packages("BiocManager")
BiocManager::install("dada2")
library("Rcpp")
library("dada2")
library("ggplot2")
```



Listing all .fastq files as FWD and REV

```
Set the path to the raw fastq data
 ```{r dataset}
path <- "/net/fs-2/scale/OrionStore/Scratch/sergroch/FON48/" # change to the directory containing the fasta files, which may
list.files(path) # To check if the fasta files in the directory
## Create lists of forward and reverse filenames for the forward and reverse reads and a list of sample names
  `{r FRList}
 Forward and reverse fasta filenames have format: SAMPLENAME_R1_001.fasta and SAMPLENAME_R2_001.fasta
LL.fnFs <- sort(list.files(path, pattern=" R1 001.fastq", full.names = TRUE))
LL.fnRs <- sort(list.files(path, pattern="_R2_001.fastq", full.names = TRUE))
# Extract sample names, assuming filenames have format: SAMPLENAME_XXX.fastq
LL.sample.names <- sapply(strsplit(basename(LL.fnFs), " "), `[`, 1)
```



Make quality scores plots

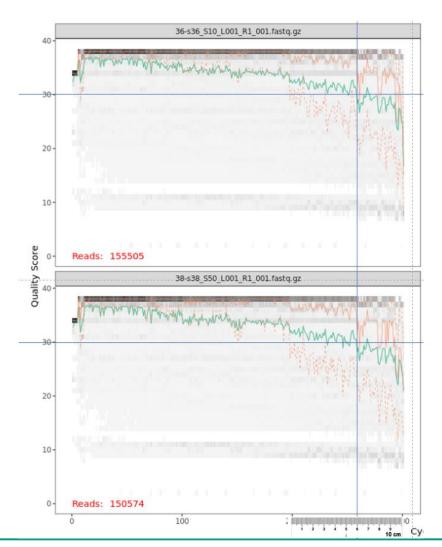
```
# This displays the forward and reverse quality scores for the samples 93 and 94
plotQualityProfile(c(LL.fnFs[93],LL.fnRs[94],LL.fnRs[94]))

To check if the foward/reverse reads from 4 samples show a consistent dropoff in quality toward the end of the reads.

```{r ReadQualityProfilesReverse}
plotQualityProfile(LL.fnFs[18:21]) #change to Rev
```

# Norwegian University of Life Sciences

#### Make quality scores plots



**Green line**: mean reads **Solid orange line**: median

**Dashed orange lines**: interquartiles

**The grey shading** is a heatmap of the frequency of the quality scores at a given position along the read. Darker shading indicates higher frequency.

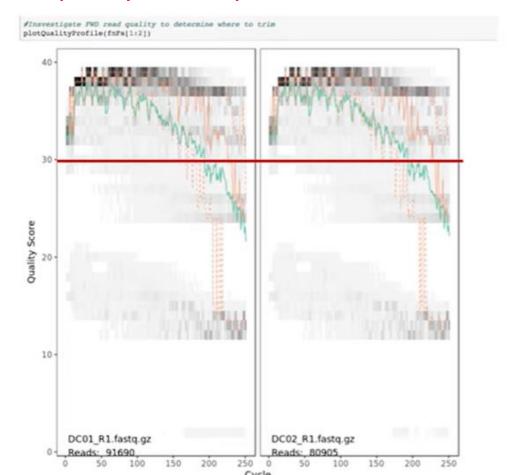
There are a few low quality scores that bring down the mean (

Oddly, there are some short sequences that are of low quality that bring down the mean at the beginning of the reads.

Note that the quality scores change along the length of the read as well as for the forward vs. the reverse reads.

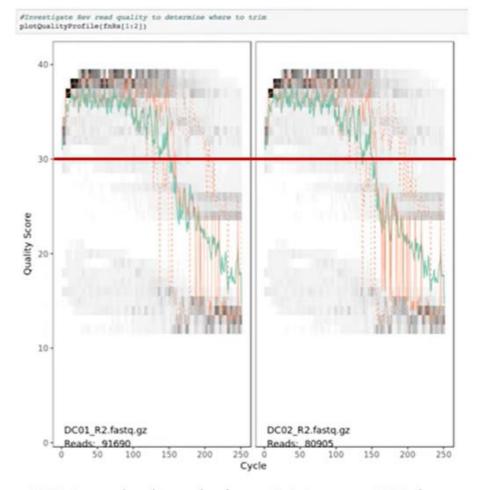
The reverse reads are normally lower quality than the forward reads.

#### Make quality scores plots



FWD reads drop below Q30 near 200 bp

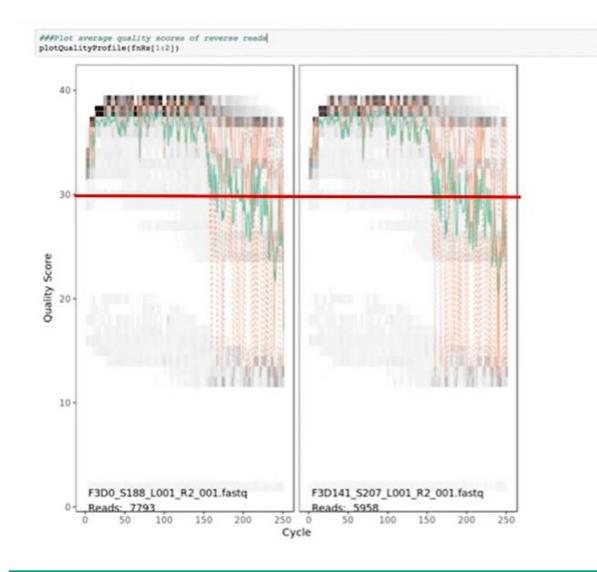




REV reads drop below Q30 near 140 bp

#### **Bioinformatics**





- The reverse reads are a different story
- This is very common reverse reads are almost always lower quality
- Overall the scores are pretty decent but they drop off right around 160 bp
- Since we'll have almost complete overlap, we can be aggressive with trimming
  - Cut at 160 bp



Create filter file path and remane file

```
Create a list of the filenames (including the path) for the filtered reads

This filename list will be later used by the filterAndTrim function from DADA2 to write the filtered reads into a new directory called "filtered".

** ** **

{r FilteredFilenames}

Place filtered files in filtered/ subdirectory

LL.filtFs <- file.path(path, "filtered", paste0(LL.sample.names, "_F_filt.fastq.gz"))

LL.filtRs <- file.path(path, "filtered", paste0(LL.sample.names, "_R_filt.fastq.gz"))
```

## Norwegian University of Life Sciences

#### Filter and Trim the reads

The **filterAndTrim** function of DADA2 is the quality control step of the pipeline.

It takes as input the forward and reverse reads, and writes a reduced set of sequences in FASTQ format that excludes low quality sequences which do not meet the criteria specified by the function arguments, **trunclen=c(275,215)**.

The **maxEE** sets the maximum expected errors in a sequence to 4 for both the forward and reverse sequences, and **rm.phix** removes any phiX DNA sequences.

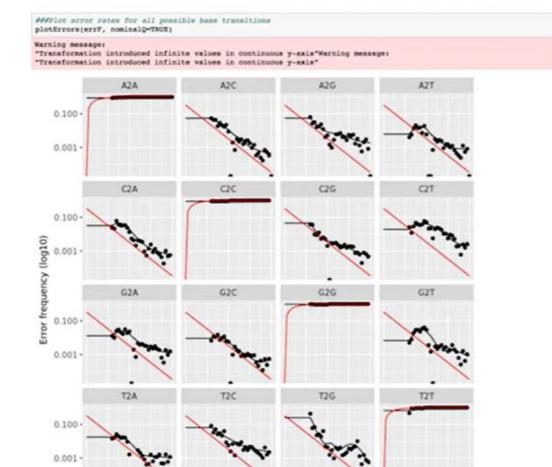
The **trimLeft** argument trims 17 bases corresponding to the length of the Pro341f primer from the beginning of the forward reads, and trims the 21 bases of Pro805r reverse primer from the beginning of the reverse reads.



#### Error rates and ploting

```
Learn the error rates for forward and reverse reads (30min)
This step applies machine learning to develop a model of the error rate for forward and reverse reads.
```{r ErrorTraining}
LL.errF <- learnErrors(LL.filtFs, multithread=FALSE)</pre>
LL.errR <- learnErrors(LL.filtRs, multithread=FALSE)
## Plot the estimated error rates for the transition types
The nucleotide transition error rates will be unique to each run.
   {r PlotErrorRates}
plotErrors(LL.errF.rds, nominalQ=TRUE)
plotErrors(LL.errR.rds, nominalQ=TRUE)
```

Error rates and ploting





- Plots frequency of each possible base transition as a function of quality score
- Black line observed error rates
- Red line expected error rate under the nominal definition of the Q-value
- In general, frequency of errors decrease as quality score increases
- These look as expected so we can proceed with data analysis

Error rates and ploting

```
Norwegian University of Life Sciences
```

```
###Estimate the error model for DADA2 algorithim using reverse reads
errR <- learnErrors(filtRs, multithread=TRUE)
Initializing error rates to maximum possible estimate.
Sample 1 - 7113 reads in 1660 unique seguences.
Sample 2 - 5463 reads in 1335 unique sequences.
Sample 3 - 2914 reads in 853 unique sequences.
Sample 4 - 2941 reads in 880 unique sequences.
Sample 5 - 4312 reads in 1286 unique sequences.
Sample 6 - 6741 reads in 1803 unique seguences.
Sample 7 - 4560 reads in 1265 unique sequences.
Sample 8 - 15637 reads in 3414 unique seguences.
Sample 9 - 11413 reads in 2522 unique sequences.
Sample 10 - 12017 reads in 2771 unique sequences.
Sample 11 - 5032 reads in 1415 unique sequences.
Sample 12 - 5299 reads in 1349 unique sequences.
Sample 13 - 18075 reads in 3290 unique sequences.
Sample 14 - 6250 reads in 1390 unique sequences.
Sample 15 - 4052 reads in 1134 unique sequences.
Sample 16 - 7369 reads in 1635 unique sequences.
Sample 17 - 4765 reads in 1084 unique sequences.
Sample 18 - 4871 reads in 1161 unique sequences.
Sample 19 - 6504 reads in 1502 unique sequences.
   selfConsist step 2
                                           Notice that error
  selfConsist step 3
  selfConsist step 4
                                           rate estimation
   selfConsist step 5
                                           takes longer for
   selfConsist step 6
                                           reverse reads
```

- This command creates an error model that will be used by the DADA2 algorithm
- Every batch of sequencing will have a different error rate
- Algorithm starts with the assumption that the error rates are the maximum possible
- Alternates error rate estimation and sample composition inference until they converge at a consistent solution

Convergence after 6 rounds. Total reads used: 135328



Dereplication: remove identical sequences

```
## Dereplication
There will be many sequences which are 100% identical. By identifying the sequences which are identical, a process called
"dereplication", the size of the dataset may be reduced.
 ``{r Dereplication}
 Dereplicate the forward and then the reverse reads
LL.derepFs <- derepFastq(LL.filtFs, verbose=TRUE)</pre>
saveRDS(LL.derepFs, "LL.derepFs.rds")
LL.derepRs <- derepFastq(LL.filtRs, verbose=TRUE)</pre>
saveRDS(LL.derepRs, "LL.derepRs.rds")
 Name the derep-class objects by the sample names
names(LL.derepFs) <- LL.sample.names
names(LL.derepRs) <- LL.sample.names
```



Sample interface (very long process) + merging

```
## Sample inference (6h)
At this step we run the principle function of the DADA2 pipeline. The dada function uses the error model along with the set
of dereplicated sequences in order to identify which sequences are likely to be real biological sequences and which are
likely the result of base-calling errors. Those sequences that are base-calling errors are clustered with the real
biological sequence from which they are likely to have derived.
```{r SampleInference}
LL.dadaFs <- dada(LL.derepFs, err=LL.errF.rds, multithread=FALSE)
saveRDS(LL.dadaFs, "LL.dadaFs.rds")
LL.dadaRs <- dada(LL.derepRs, err=LL.errR, multithread=FALSE)</pre>
saveRDS(LL.dadaRs, "LL.dadaRs.rds")
LL.dadaFs[[1]]
Merge paired reads (10min)
 ```{r MergeReads}
LL.mergers <- mergePairs(LL.dadaFs, LL.derepFs, LL.dadaRs, LL.derepRs, verbose=TRUE)
head(LL.mergers[[79]])
```



Sample interface (very long process) + merging

```
###Merge denoised reads
mergers <- mergePairs(dadaFs, derepFs, dadaRs, derepRs, verbose=TRUE)
# Inspect the merger data.frame from the first sample
head(mergers[[1]])

6594 paired-reads (in 104 unique pairings) successfully merged out of 7113 (in 254 pairings) input.
5047 paired-reads (in 78 unique pairings) successfully merged out of 5463 (in 199 pairings) input.
2663 paired-reads (in 52 unique pairings) successfully merged out of 2914 (in 142 pairings) input.
2574 paired-reads (in 54 unique pairings) successfully merged out of 2941 (in 163 pairings) input.
3668 paired-reads (in 53 unique pairings) successfully merged out of 4312 (in 203 pairings) input.</pre>
```

- Merges paired end reads only if they exactly overlap
 - This is because both forward and reverse reads have been denoised and should be error-free
 - Can be changed by adding maxMismatch option
- By default, the program requires 20 nt of overlap but you can lower it if need be with minOverlap

Norwegian University of Life Sciences

27

Amplicon Sequence Variant table

```
## Make the Amplicon Sequence Variant (ASV) table
DADA2 produces Amplicon Sequence Variants (ASVs). Hence, here we produce an ASV table.
LL.seqtab <- makeSequenceTable(LL.mergers)</pre>
# This produces the dimensions of the table, with the rows as samples and columns as ASVs.
dim(LL.seqtab)
table(nchar(getSequences(LL.seqtab)))
rownames(LL.seqtab)[1]
colnames(LL.seqtab)[1]
```



Remove quimeras

```
## Remove chimeras from the ASV table (1h)

Chimeras occur when primer extentions truncated during a PCR cycle anneal to mismatched 16S template from another taxon and prime the extension of the truncated fragment to create a chimeric sequence consisting of partial 16S sequence from 1 taxon attached to 16S from a different taxon. These chimeric sequences thus do not represent real biological sequences and must be removed.

**TRemoveChimeras**
LLSeqtab.nochim <- removeBimeraDenovo(LL.seqtab, method="consensus", multithread=FALSE, verbose=TRUE)

# This Lists the dimensions of the ASV table following chimera removal. It is the number of samples X number of ASVs.

dim(LLSeqtab.nochim)

# The sum function sums the sequence counts from each combination of ASV and sample. Thus, it gives the number of sequences in each table.

sum(LLSeqtab.nochim)/sum(LL.seqtab)

saveRDS(LLSeqtab.nochim, "LLSeqtab.nochim.rds")
```



Notes for Chimera Checking

- It is normal for a large majority of <u>unique sequences</u> to be flagged as chimeric
 - They should make up a small proportion of total reads
- If a majority of total reads are chimeric, you have a problem
 - Most likely explanation primers were not completely removed from reads
 - Using primers with ambiguous bases can cause reads to be flagged as chimeras
 - Re-run filterAndTrim() command using the trimLeft argument
 - If that doesn't fix it, there may be other factors at play
 - Try trimming more low quality bases
 - Think about which hypervariable region you're sequencing (i.e. V4 vs V4V5)
 - How complex is your community?



Assigning Taxonomy

```
## Assign taxonomy (2.5h)
This step determines the taxonomy of each ASV. Different reference databases may be used, and there are several formatted
for use with DADA2 and available at
<a href="https://benjjneb.github.io/dada2/training.html">https://benjjneb.github.io/dada2/training.html</a>). Here we use the
latest version of the Silva database. Note that there are entries in Silva which do not conform to 7 hierarchy levels. Also
note that identifical down to the species level requires a separate reference file and algorithm.
   {r AddTaxonomy}
LL.taxa <- assignTaxonomy(LLSeqtab.nochim, paste0("/net/fs-2/scale/OrionStore/Home/sergroch/16S
material/silva nr99 v138.1 wSpecies train set.fa.gz"), multithread=FALSE)
# A separate approach is needed for assigning species designations #silva species assignment v138.1.fa.gz
LL.taxa <- addSpecies(LL.taxa, paste0("/net/fs-2/scale/OrionStore/Home/sergroch/16S
material/silva species assignment v138.1.fa.gz"))
taxa.print <- LL.taxa # Removing sequence rowngmes for display only
rownames(taxa.print) <- NULL
head(taxa.print)
saveRDS(LLSeqtab.nochim, "Seqtab.nochim.rds")
#Seatab.nochim <- readRDS("SergioSeatab.nochim.rds")
saveRDS(LL.taxa, "LL.taxa.rds")
```

Norwegian University of Life Sciences

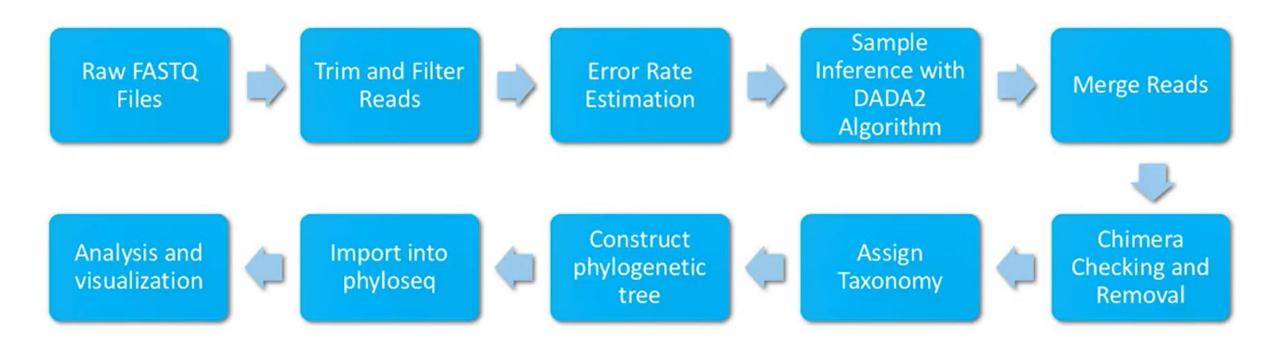
Taxonomy table

\square	Α	В	С	D	Е	F	G	Н	J	K	L	М	N	0	Р
1	Row.names	Kingdom	Phylum	Class	Order	Family	Genus	Species	S1	S10	S11	S12	S13	S14	S15
2	TAGGGAATC	Bacteria	Firmicutes	Bacilli	Lactobacillal	es			C	3224	2715	4204	39484	33014	21429
3	TGGGGAATA	Bacteria	Actinobacter	Actinobacter	Micrococcale	Beutenbergi	Salana		C	347	0	400	3900	2976	3608
4	TGGGGAATA	Bacteria	Proteobacter	Gammaprote	Xanthomona	Xanthomona	Silanimonas		C	0	0	0	0	0	0
5	TAGGGAATC	Bacteria	Firmicutes	Bacilli	Lactobacillal	Carnobacter	Granulicatell	elegans	C	0	0	0	32	0	10
6	TGGGGAATA	Bacteria	Proteobacter	Alphaproteo	Sphingomon	Sphingomon	Sphingomon	as	C	0	0	0	0	0	0
7	TGGGGAATT	Bacteria	Proteobacter	Gammaprote	Burkholderia	Comamonad	Comamonas		C	0	0	0	0	0	0
8	TGGGGAATT	Bacteria	Proteobacter	Gammaprote	Burkholderia	Comamonad	Acidovorax	facilis	C	0	0	0	0	0	0
9	TGAGGAATA	Bacteria	Bacteroidota	Bacteroidia	Flavobacteri	Flavobacteri	Myroides	injenensis	C	0	0	0	0	33	0
10	TGGGGAATA	Bacteria	Firmicutes	Clostridia	Clostridiales	Clostridiacea	Hathewaya		C	0	0	0	0	0	0
11	TGAGGAATA	Bacteria	Bacteroidota	Bacteroidia	Bacteroidale	Paludibacter	aceae		C	0	0	0	0	0	0
12	TGGGGAATA	Bacteria	Actinobacter	Actinobacter	Micrococcale	Micrococcac	Glutamicibad	creatinolytic	C	0	0	0	0	0	0

Bioinformatics



DADA2 Workflow



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

Norwegian University of Life Sciences

Packages needed

```
## Load the required package. Install ggstatsplot before running the script. If you are working with orion, this should be done with R.4.3.1 as this was recently updated for installing the package
```{r loading package, echo=FALSE}
library(Rcpp)
library(dada2)
library(permute)
library(permute)
library(vegan)
library(ggplot2)
library(tidyverse)
library(ggstatsplot)
library(dplyr)
library(openxlsx)
```



#### Create a phyloseq object

A phyloseq object consists of the 5 main data types needed for complete microbiome analysis.

- An ASV table like the one produced
- The sample metadata table containing, for example, fish number, tank, diet, [DNA],.....
- A reference nucleotide sequence for each ASV
- A phylogenetic tree
- A taxonomy table with the levels of the taxonomic hierarchy for every ASV



#### Create a phyloseq object

```
##Create a phyloseq object
 {r MakePSObject}
ps <- phyloseq(otu table(count tab, taxa are rows=FALSE),</pre>
 sample data(LL.samdf),
 tax table(tax tab))
##Remove the undetermined sample & other if necessary
ps <- prune samples(sample names(ps) != "Undetermined F filt.fastq.gz", ps) # Remove undetermined sample
##Transform the phyloseq object into proportion
ps_tss <- transform_sample_counts(ps, function(x){x / sum(x)})</pre>
```



#### Create a phyloseq object

```
Taxonomy-based filtering
Remove features without a phylum-level annotation and those assigned as chloroplast or mitochondria. Note that the taxonomic
labels are database specific and may change in different versions of the same database. Make sure you're using the correct
taxonomic labels to remove chloroplast and mitochondria.
ps tss <- subset taxa(ps tss, !is.na(Phylum) & !Phylum %in% c("", "unassigned")) %>%
 subset taxa(Order != "Chloroplast" | is.na(Order)) %>%
 subset taxa(Family != "Mitochondria" | is.na(Family))
Prevalence-based filtering
Features that show up in only one or a few samples may not represent real biological diversity but rather PCR/sequencing
errors (such as PCR chimeras) or reagent contaminants.
ps_tss <- subset_samples(ps_tss, !Sample kind %in% c("negative control")) %>%
 filter taxa(., function(x) sum(x > 0) > 1, TRUE) %>%
 taxa names() %>%
 prune taxa(ps tss)
phyloseq-class experiment-level object
otu table() OTU Table: [1617 taxa and 96 samples]
sample data() Sample Data: [96 samples by 8 sample variables]
tax table() Taxonomy Table:
 [1617 taxa by 8 taxonomic ranks]
```

# Norwegian University of Life Sciences

#### Filter contaminants

```
Filter contaminants
Screening of reagent contaminants
The screening of reagent contaminants will be based on two typical characteristics of contaminanting sequences as outlined
in the paper [Simple statistical identification and removal of contaminant sequences in marker-gene and metagenomics
data | (https://microbiomejournal.biomedcentral.com/articles/10.1186/s40168-018-0605-2): they are likely to have frequencies
that inversely correlate with sample DNA concentration and are likely to have higher prevalence in control samples than in
true samples. The authors developed an R package, [*decontam*](https://github.com/benjjneb/decontam), for removing
contaminating sequences in the marker-gene and shotgun metagenomics data. The package, however, does not make use of
positive controls for the identification of contaminating sequences. As removing of features may critically affect
downstream analyses, we'll do it by manual screening based on the aforementioned principles.
Identify reagent contaminants
Data wrangling
Make a dataframe containing features present in the negative controls and mock samples.
decontam <- ps tss %>%
 subset samples(Diet %in% c("FP", "Mock", "NC")) %>%
 filter taxa(., function(x) sum(x > 0) > 0, TRUE) %>%
 taxa names() %>%
 prune taxa(ps tss) %>%
```

**FOODS PNORWAY** 

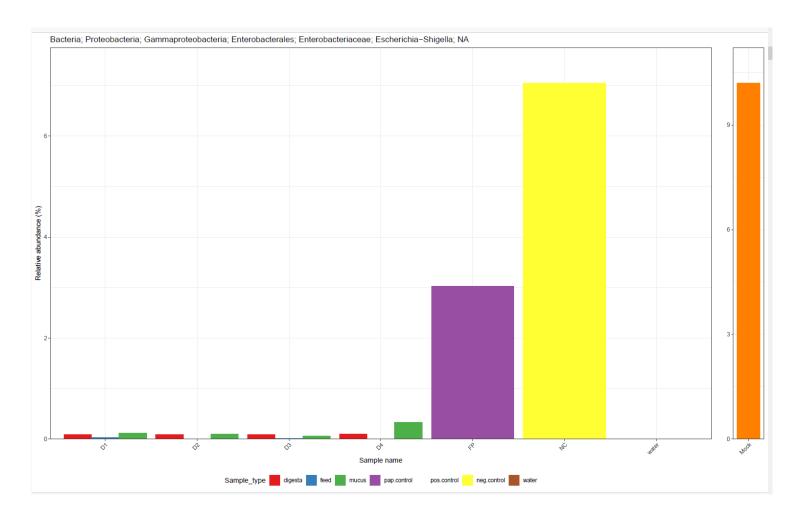


### Filter contaminants (prevalence)

```
Prevalence-based classification - this generates a graphic of pdf file. Here we use barplots to visualize the abundance
and prevalence of the features found in the control samples After the production of the bar plots, manually inspect them to
identify the contaminants in the negative control.
 {r, results='hide'}
library(cowplot)
decontam spl1 <- group split(decontam, OTU)</pre>
pdf("prevalence contam.pdf", width = 16, height = 10)
lapply(seq_along(decontam_spl1), function(x){
 p1 <- filter(decontam spl1[[x]], Sample_type != "pos.control") %>%
 plot prevalence(x = Diet, y = Abundance, bar color = Sample type,
 xlab = "Sample name", ylab = "Relative abundance (%)",
 title = unique(decontam spl1[[x]][, "tax"]))
 # make a bar plot with mock only
 p2 <- filter(decontam_spl1[[x]], Sample_type == "pos.control") %>%
 plot prevalence(x = Diet, y = Abundance, bar color = Sample type, xlab = "", ylab = "")
 plot grid(p1, p2 + theme(legend.position = "none"), nrow = 1, align = 'h', axis = "bt", rel widths = c(13, 1))
dev.off()
```

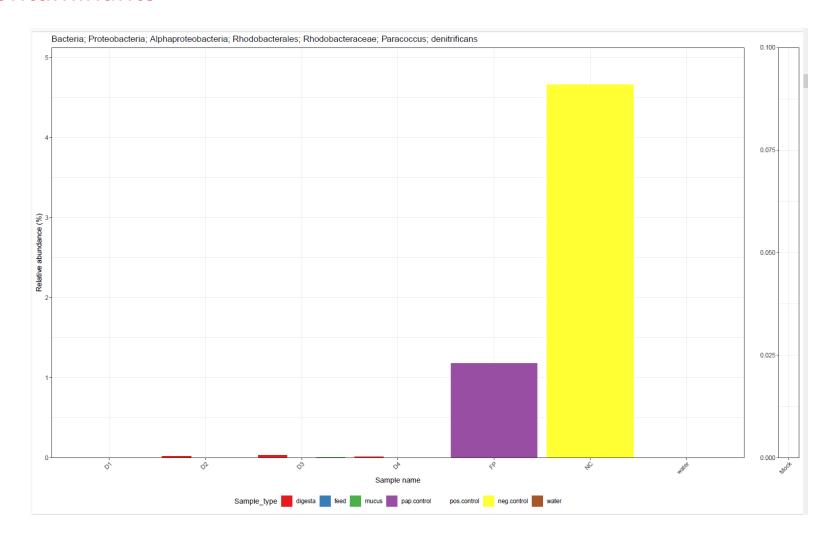
# Norwegian University of Life Sciences

### Filtration of contaminants



# Norwegian University of Life Sciences

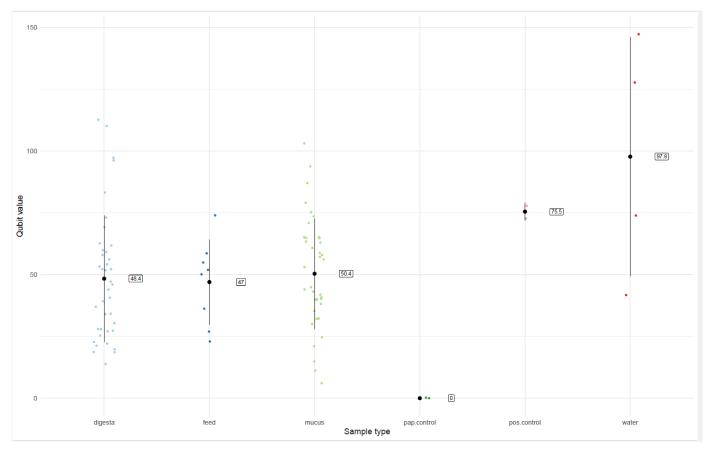
### Filtration of contaminants





### Filter contaminants (DNA concentration)

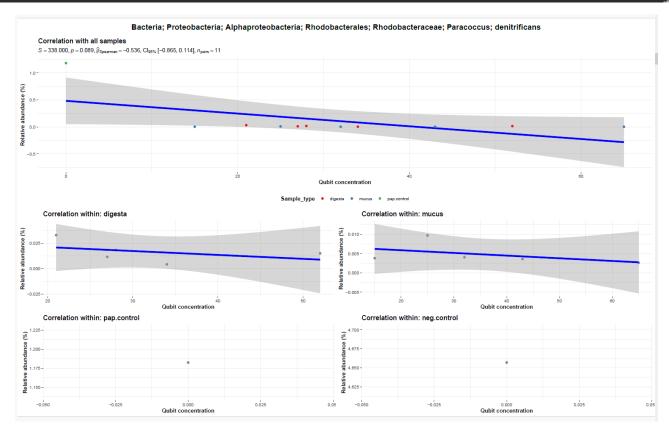
## Inspect bacterial DNA concentration ##Before we proceed with the identification of contaminating features, let's look at the DNA concentration as measured by Qubit values of the DNA templates used for the amplicon PCR.





### Filter contaminants (frequency)

### Frequency-based classification
##Here we visualize correlations between the DNA concentration values and the relative abundance of features found in the
control samples. Features showing inverse correlations with DNA concentration are potential contaminating features
introduced during the amplicon PCR, which is the main source of reagent contamiantion in this study.



FOODS PNORWAY 42

# Norwegian University of Life Sciences

#### Determine and remove contaminants

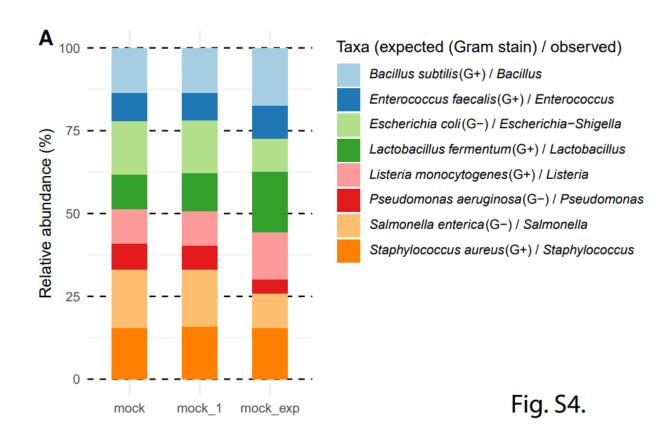


#### Determine and remove contaminants

```
##Check the distibution of contaminating features.
 `{r, fig.width=16, fig.height=10}
prune taxa(taxa names(ps tss) %in% contam mock$OTU, ps tss) %>%
 plot bar(x = "Sample type", fill = "Phylum", title = "Cross-contaminating features in the mock samples") +
 scale y continuous(expand = expansion(mult = c(0, 0.05))) +
 scale fill brewer(palette = "Paired") +
 theme bw() +
 theme(legend.position = "bottom", axis.text.x = element text(angle = 90, hjust = 1))
Remove reagent and cross-contaminants
ps_nocontam <- prune_taxa(taxa_names(ps_tss_nocontam), ps)</pre>
ps nocontam mock <- subset samples(ps nocontam, Sample type == "pos.control")</pre>
ps nocontam mock <- prune taxa(!taxa names(ps nocontam mock) %in% contam mock$OTU, ps nocontam mock)
ps nocontam <- subset samples(ps nocontam, Sample type != "pos.control") %>%
 merge loseq(ps nocontam mock)
```

Norwegian University
of Life Sciences

Positive control: Mock





### More packages to install

library(devtools)

## Load the required package. Install ggstatsplot before running the script. If you are working with orion, this should be done with R.4.0.4 as this was recently updated for installing the package ``{r loading package, echo=FALSE} library(Rcpp) library(dada2) .ibrary(phyloseq) ibrary(permute) ibrary(lattice) .ibrary(vegan) .ibrary(ggplot2) library(tidyverse) library(ggstatsplot)library(dplyr) library(microbiome) ibrary(microbiomeutilities) library(knitr) library(RColorBrewer) library(DT) library(cowplot) library(PerformanceAnalytics) library(venn) library(philr) library(MicrobeR)

FOODS PNORWAY 46



### Extract ps\_nocontam and check coverage

```
##extract feature table, taxonomy and metadata from the phyloseq object. Check if your feature table (count_tab) are sample
ID as rows. check with taxa are rows(phyloseq object). if the answer is FALSE, transpose it as done below. otherwise, don't
need to do it.
count tab <- as.data.frame(t(otu table(ps nocontam))) #to transpose if the observation do not mach to run tab phy
#count tab <- as.data.frame((otu table(ps nocontam))
tax tab <- tax table(ps nocontam) %>% as("matrix") %>% as.data.frame()
metadata <- data.frame(sample data(ps nocontam), check.names = FALSE)</pre>
Overview of taxonomy assignments
First of all, let's look at the coverage of taxonomy assignments at different levels.
library(tidyr)
tax tab %>%
 gather("Rank", "Name", rank names(ps nocontam)) %>%
 group by(Rank) %>%
 # Empty taxonomic ranks may be na or strings containing "uncultured" or "Ambiguous_taxa"
 summarize(ASVs classified = sum(!is.na(Name) & !grepl("uncultured|Ambiguous|metagenome", Name))) %>%
 mutate(Frac classified = ASVs classified / ntaxa(ps nocontam),
 Frac classified = ifelse(Frac classified == 1, "100", round(Frac classified * 100, 1)),
 Frac classified = paste(Frac classified, "%"),
 Rank = factor(Rank, rank names(ps nocontam))) %>%
 arrange(Rank) %>%
 datatable(options = list(columnDefs = list(list(className = 'dt-left', targets = c(0:3)))))
```



Extract ps\_nocontam and check coverage

Show	10 v entries		Search:	Search:			
	Rank	ASVs_classified	Frac_classified				
1	Kingdom	1617	100 %				
2	Phylum	1617	100 %				
3	Class	1584	98 %				
4	Order	1558	96.4 %				
5	Family	1495	92.5 %				
6	Genus	1311	81.1 %				
7	Species	203	12.6 %				
Showin	ng 1 to 7 of 7 entries	Previous 1	Next				

**FOODS PNORWAY** 



Group samples by type (phylum and genus)

```
row.names(count tab)
openxlsx::write.xlsx(tab_phy, file = "phylum.xlsx", overwrite = TRUE)
tab_phy1d <- tab_phy %>%
 select(c(01,02,03,04,05,06,07,08,09,10,11,12,13,14,15,16,17,18,19,20,21,22,23,24,25,26,27,28,29,30,31,33,34,35,
 37,46,47,48)) # removed sample 36, fish 10
tab phy1m <- tab phy %>%
 select(c(38,39,40,41,42,44,45,49,50,51,52,53,54,55,57,58,59,60,61,62,63,64,65,66,67,68,69,70,71,72,73,74,75,76,
77,78,79,80)) # removed sample 43, fish 10
tab phy2 <- tab phy %>%
 select(c(81,82,83,84,85,86,87,88))
tab_phy3 <- tab_phy %>%
 select(c(89,90,91,92))
openxlsx::write.xlsx(tab phy1d, file = "phylum digesta.xlsx",overwrite = TRUE)
openxlsx::write.xlsx(tab phy1m, file = "phy1um mucus.xlsx",overwrite = TRUE)
openxlsx::write.xlsx(tab phy2, file = "phylum feed.xlsx",overwrite = TRUE)
openxlsx::write.xlsx(tab_phy3, file = "phylum_water.xlsx",overwrite = TRUE)
tab phya <- tab phy %>% rownames to column()
openxlsx::write.xlsx(tab_phya, file = "phylum_name.xlsx",overwrite = TRUE)
```

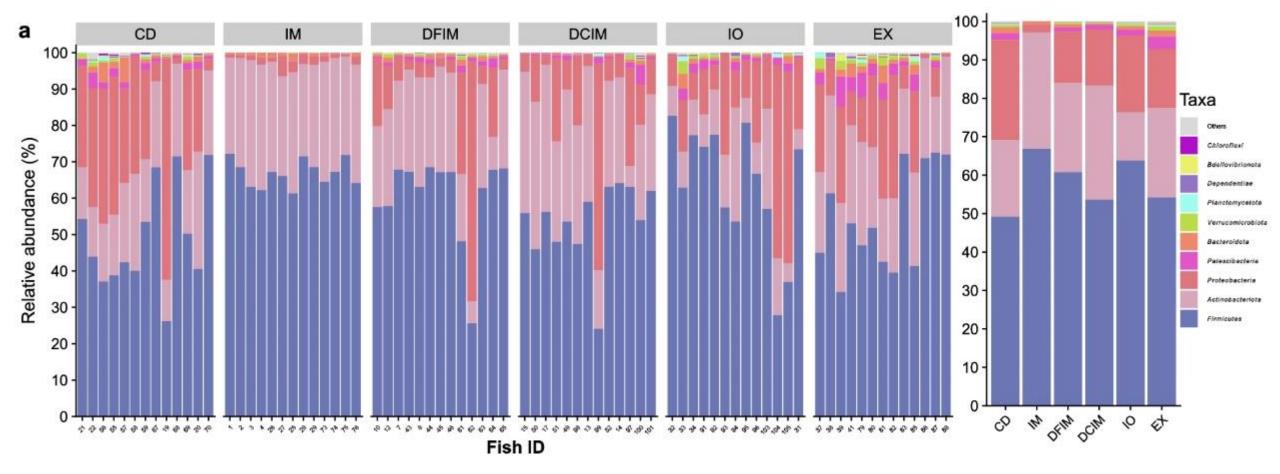


```
##Make taxa barplot at phylum level - DIGESTA
metadata1d <- subset (metadata, Sample type == "digesta")</pre>
metadata1d$Diet <- factor(metadata1d$Diet, levels = c("D1", "D2", "D3", "D4"))</pre>
p phy ind <- make taxa barplot(table = tab phy1d,</pre>
 metadata = metadata1d,
 group_by = factor(Diet, c("D1", "D2", "D3", "D4")),
 ntaxa = 10,
 nrow = 1,
 plot mean = FALSE,
 cluster sample = FALSE,
 sample \overline{label} = fish,
 italize taxa name = FALSE,
 colors=colors)
p phy ind <- p phy ind + labs(x = "fish")
p phy meand <- make taxa barplot(table = tab phy1d,
 metadata = metadata1d,
 group_by = Diet,
 ntaxa = 10,
 nrow = 1,
 plot mean = TRUE,
 cluster sample = FALSE,
 italize taxa name = TRUE,
 colors=colors)
p phy meand \leftarrow p phy meand + labs(x = "", y = "")
```

FOODS PNORWAY 50

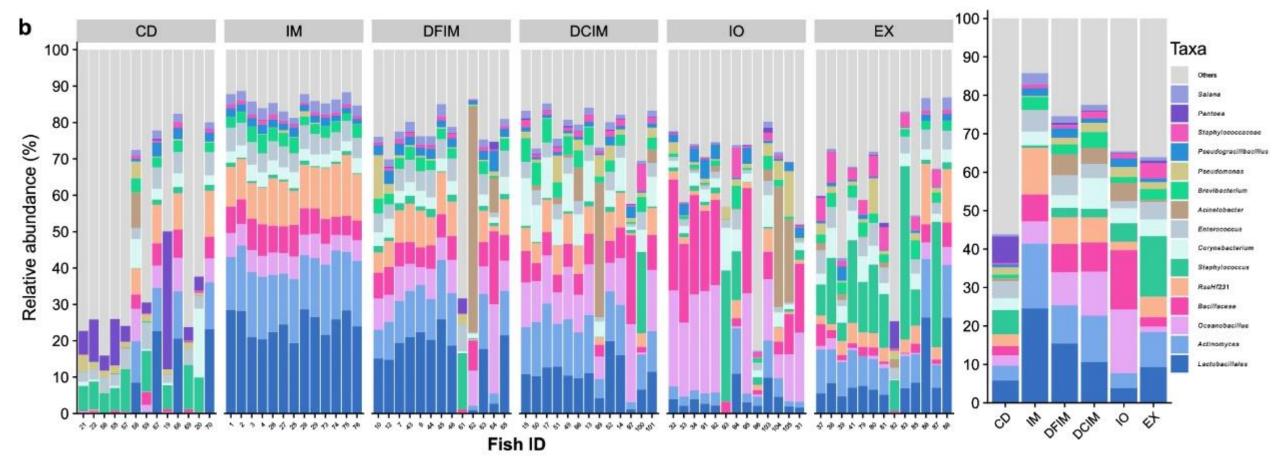


Digesta - Phylum



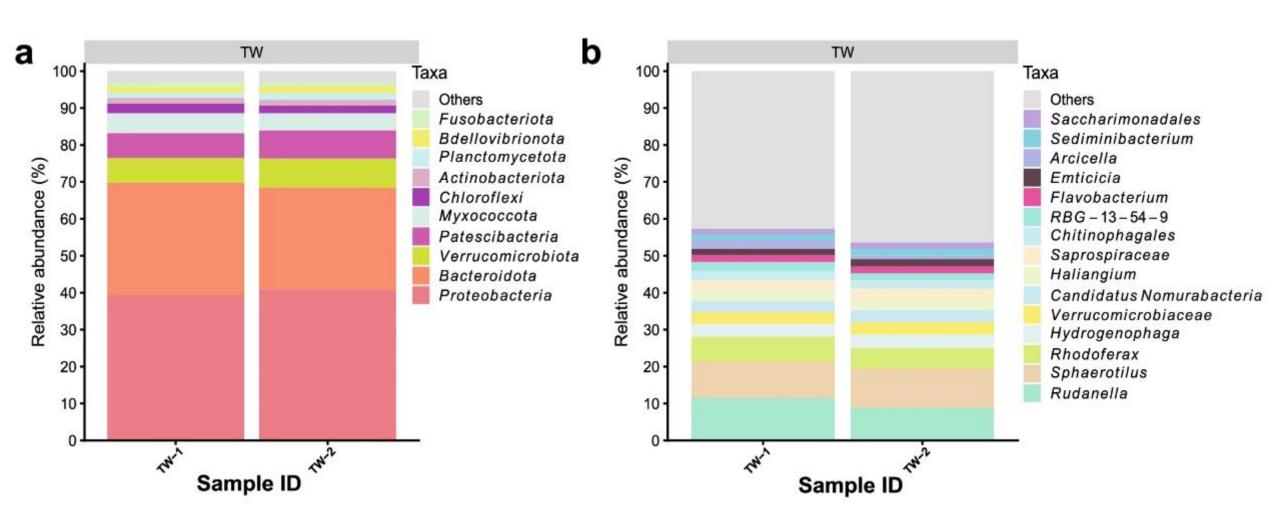


Digesta - Genus



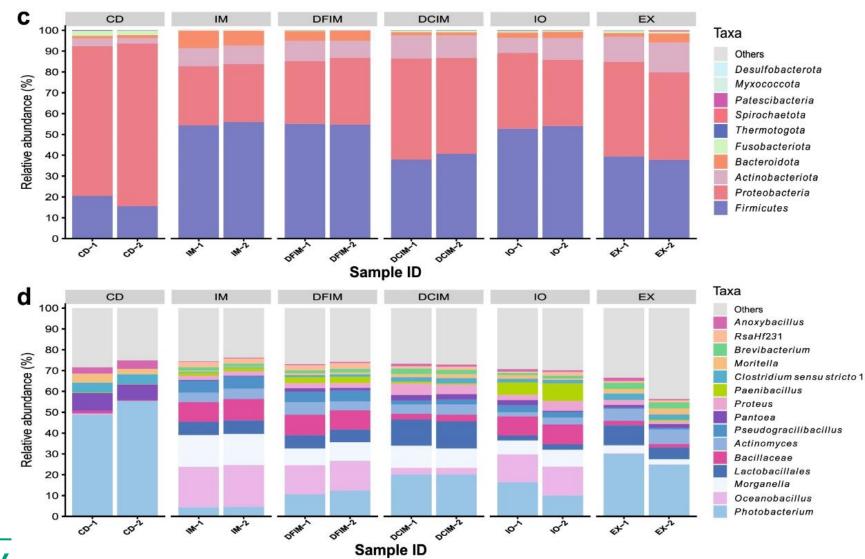
# Norwegian University of Life Sciences

Water samles



# Norwegian University of Life Sciences

### Feed samles





#### Core Microbiome

#### Calculate feature prevalence

```
Compute feature prevalence
CD <- subset_samples(ps_nocontam, Diet == "CD") %>% prevalence()
IM <- subset_samples(ps_nocontam, Diet == "IM") %>% prevalence()
DFIM <- subset_samples(ps_nocontam, Diet == "DFIM") %>% prevalence()
DCIM <- subset_samples(ps_nocontam, Diet == "DCIM") %>% prevalence()
IO <- subset_samples(ps_nocontam, Diet == "IO") %>% prevalence()
EX <- subset_samples(ps_nocontam, Diet == "EX") %>% prevalence()
```

##Get core features that are present in at least 80% samples under the different diets

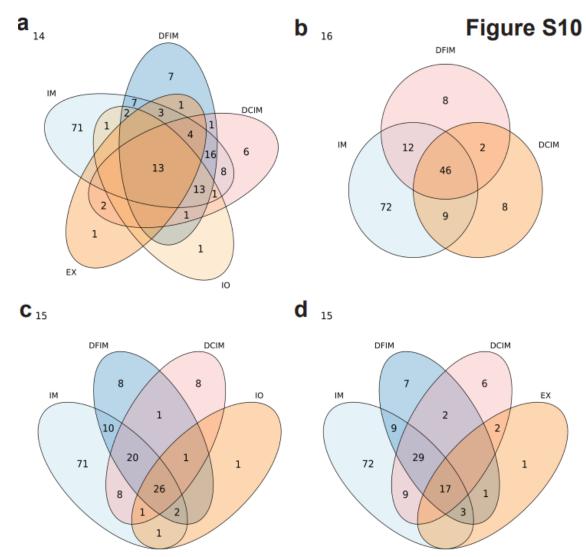
```
core_taxa <- cbind.data.frame(CD, IM, DFIM, DCIM, IO, EX) %>%
 rownames_to_column("featureID") %>%
get core features based on 80% prevalence threshold
 filter(CD >= 0.8|IM >= 0.8|DFIM >= 0.8|DCIM >= 0.8|IO >= 0.8|EX >= 0.8)
```

##Add taxonomy to core features

```
core_taxa_tab <- rownames_to_column(tax_tab, "featureID") %>%
 inner_join(core_taxa, by = "featureID") %>%
 mutate(prev_all = rowSums(.[9:14])) %>%
 arrange(desc(prev_all)) %>%
 mutate if(is numeric _~ifelse(v == 1 _100 _round(v * 100 _1))) %>%
```



Core Microbiome



Digesta, Feed and water overlap

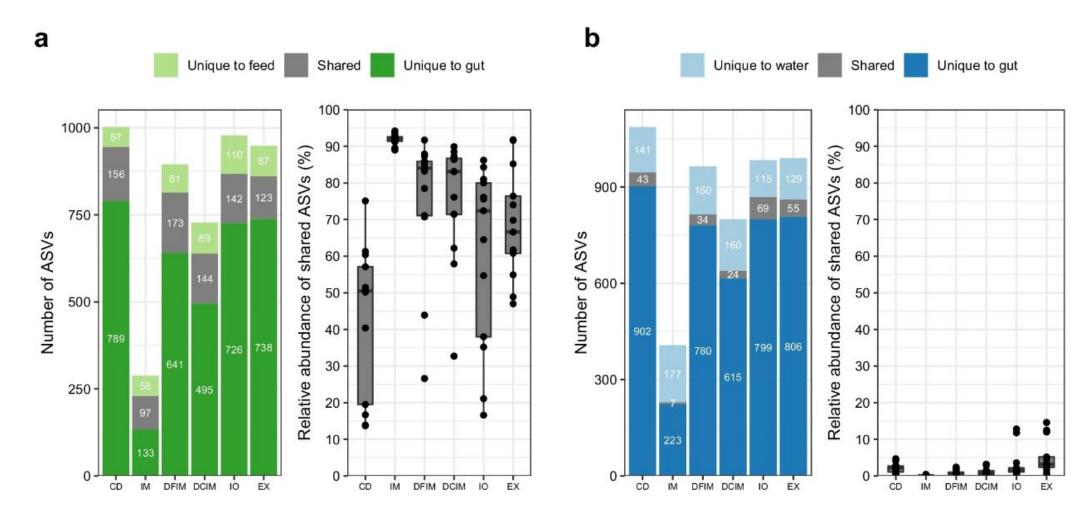


#### Compute ASV overlap

```
ovrl <- bind_cols(prvl) %>%
 pivot_longer("CD":"EX",
 names_to = "Sample",
 values_to = "Gut") %>%
 # compute ASV overlap between water/feed and intestinal samples
 mutate(
 WaterOverlap = case_when(
 Water > 0 & Gut > 0 ~ "Shared",
 Water > 0 & Gut == 0 ~ "Unique to water",
 Water == 0 & Gut > 0 ~ "Unique to gut",
 TRUE ~ "Absent"),
 FeedOverlapCD = case_when(
 FeedCD > 0 & Gut > 0 ~ "Shared",
 FeedCD > 0 & Gut == 0 ~ "Unique to feed",
 FeedCD == 0 & Gut > 0 ~ "Unique to gut",
 TRUE ~ "Absent"),
 FeedOverlapIM = case_when(
 FeedIM > 0 & Gut > 0 ~ "Shared",
 FeedIM > 0 & Gut == 0 ~ "Unique to feed",
 FeedIM == 0 & Gut > 0 ~ "Unique to gut",
 TRUE ~ "Absent"),
 FeedOverlapDFIM = case_when(
 FeedDFIM > 0 & Gut > 0 ~ "Shared",
 FeedDFIM > 0 & Gut == 0 ~ "Unique to feed",
 FeedDFIM == 0 & Gut > 0 ~ "Unique to gut",
 TRUE ~ "Absent"),
 FeedOverlapDCIM = case_when(
 FeedDCIM > 0 & Gut > 0 ~ "Shared",
 FeedDCIM > 0 & Gut == 0 ~ "Unique to feed",
```

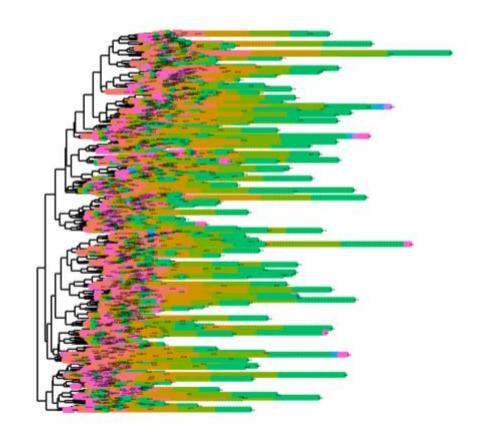


Digesta, Feed and water overlap





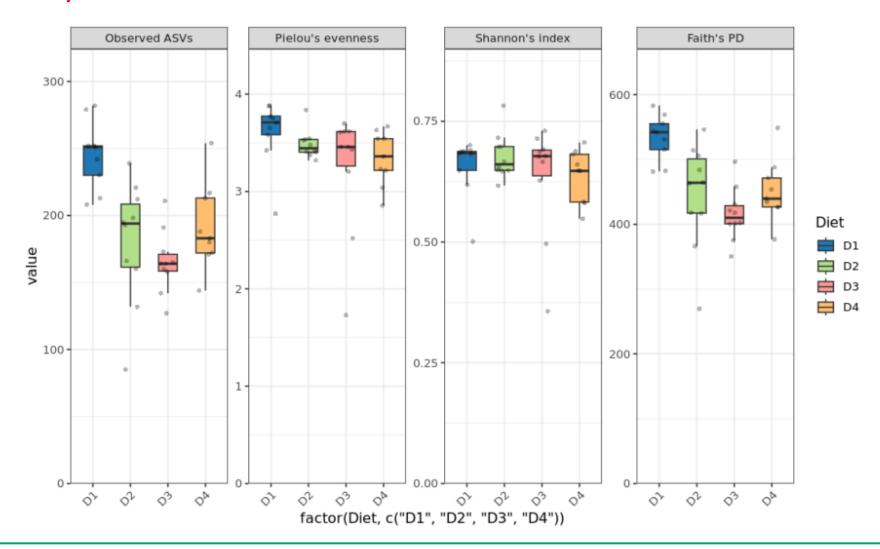
More and more packages + plot phylogenetic tree





# Norwegian University of Life Sciences

Plot alpha-diversity indices





#### statistics

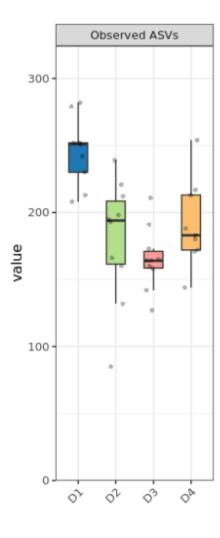
```
Alpha-diversity group significance
The statisitical difference between the diets for the four alpha diversity measurements were evaluated using kruska-walis
test and the significance between diets were identified with pairwise wilcox test.
##Flitering the adiv dataframe into each alpha diversity index
 ##Filter only the <u>dataframe</u> for Observed <u>ASVs</u>
adiv o <- adiv %>%
 filter(adiv == "Observed ASVs")
##Filter only the dataframe for Pielou's evenness
adiv p <- adiv %>%
 filter(adiv == "Pielou's evenness")
##Filter only the dataframe for Shannon's index
adiv s <- adiv %>%
 filter(adiv == "Shannon's index")
##Filter only the dataframe for Faith's PD
adiv f <- adiv %>%
 filter(adiv == "Faith's PD")
```

FOODS PNORWAY 61

#### statistics

```
Test the diet effect on alpha diversity indices - Kruska-walis test
kruskal.test(value ~ Diet, data = adiv o)
kruskal.test(value ~ Diet, data = adiv p)
kruskal.test(value ~ Diet, data = adiv s)
kruskal.test(value ~ Diet, data = adiv_f)
 Kruskal-Wallis rank sum test
data: value by Diet
Kruskal-Wallis chi-squared = 18.146, df = 3, p-value = 0.0004103
 Kruskal-Wallis rank sum test
data: value by Diet
Kruskal-Wallis chi-squared = 6.7913, df = 3, p-value = 0.07886
 Kruskal-Wallis rank sum test
data: value by Diet
Kruskal-Wallis chi-squared = 1.7168, df = 3, p-value = 0.6332
 Kruskal-Wallis rank sum test
data: value by Diet
Kruskal-Wallis chi-squared = 16.681, df = 3, p-value = 0.0008218
```



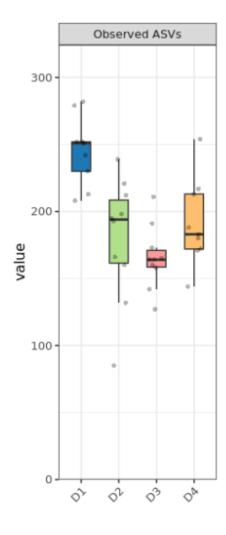


#### statistics

```
Pair-wise comparison
from rhe P-values for alpha diversity indices, observed ASVs was the only one below 0.05. Thus, we proceed to use wilcox
pairwise comparison to identify differences between diets.
setwd("/net/fs-2/scale/OrionStore/Scratch/sergroch/FON48/")
 ☆ 🗷 🕽
wilcox o <- compare means(value ~ Diet, adiv o, method = "wilcox.test")</pre>
wilcox p <-compare means(value ~ Diet, adiv p, method = "wilcox.test")</pre>
wilcox_s <-compare_means(value ~ Diet, adiv_s, method = "wilcox.test")</pre>
wilcox_f <-compare_means(value ~ Diet, adiv_f, method = "wilcox.test")</pre>
illustrator on any of the editing software later
write.csv(wilcox_o, file = "Wilcox pairwise comparison_Observed ASVs_digesta", row.names = FALSE)
write.csv(wilcox p, file = "Wilcox pairwise comparison Pielou's evenness digesta", row.names = FALSE)
write.csv(wilcox_s, file = "Wilcox pairwise comparison_Shannon's index_digesta", row.names = FALSE)
write.csv(wilcox f, file = "Wilcox pairwise comparison Faith's PD digesta", row.names = FALSE)
openxlsx::write.xlsx(wilcox_o, file = "Wilcox pairwise comparison_Observed ASVs_digesta.xlsx", overwrite = TRUE)
openxlsx::write.xlsx(wilcox_p, file = "Wilcox pairwise comparison_Pielou's evenness_digesta.xlsx", overwrite = TRUE)
openxlsx::write.xlsx(wilcox s, file = "Wilcox pairwise comparison Shannon's index_digesta.xlsx", overwrite = TRUE)
openxlsx::write.xlsx(wilcox f, file = "Wilcox pairwise comparison Faith's PD digesta.xlsx", overwrite = TRUE)
```

A tibble: 6 x 8									
<b>.y.</b> <chr></chr>	group1 <chr></chr>	group2 <chr></chr>	<b>p</b> <dbl></dbl>	p.adj <dbl></dbl>	p.format <chr></chr>	p.signif <chr></chr>	method <chr></chr>		
value	D4	D1	2.756067e-03	0.01400	0.0028	**	Wilcoxon		
value	D4	D3	6.525363e-02	0.20000	0.0653	ns	Wilcoxon		
value	D4	D2	1.000000e+00	1.00000	1.0000	ns	Wilcoxon		
value	D1	D3	8.660071e-05	0.00052	8.7e-05	***	Wilcoxon		
value	D1	D2	5.672346e-03	0.02300	0.0057	**	Wilcoxon		
value	D3	D2	1.654939e-01	0.33000	0.1655	ns	Wilcoxon		
6 rows									







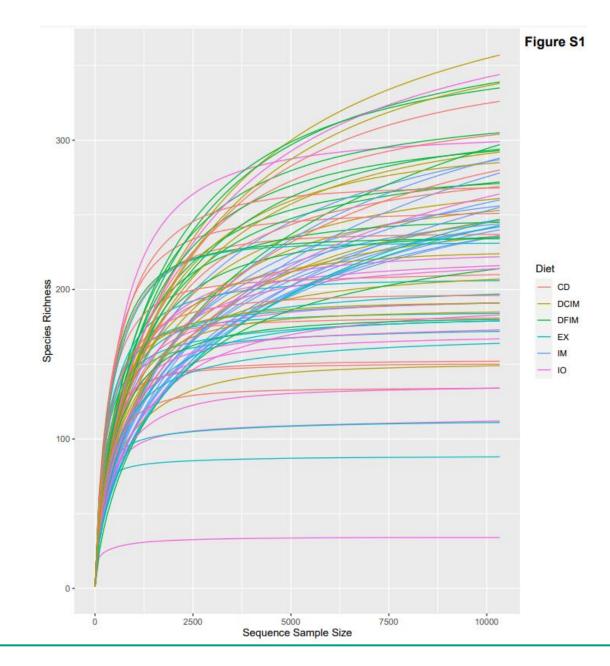
More packages again + rarefraction

```
ЭЩ
 <u>:</u> ш
Check sample rarefraction curve.
p <- ggrare(pseq, step = 1000, color ="Diet", label = "Diet", se = FALSE)</pre>
 <- p + facet wrap(~Diet)</pre>
pseq1 <- subset samples(pseq, Diet == "D1")</pre>
p1 <- ggrare(pseq1, step = 1000, color ="Diet", label = "Diet", se = FALSE)
pseq2 <- subset samples(pseq, Diet == "D2")</pre>
p2 <- ggrare(pseq2, step = 1000, color ="Diet", label = "Diet", se = FALSE)
pseq3 <- subset samples(pseq, Diet == "D3")</pre>
p3 <- ggrare(pseq3, step = 1000, color ="Diet", label = "Diet", se = FALSE)
pseq4 <- subset samples(pseq, Diet == "D4")</pre>
p4 <- ggrare(pseq4, step = 1000, color ="Diet", label = "Diet", se = FALSE)
 Library Size
```

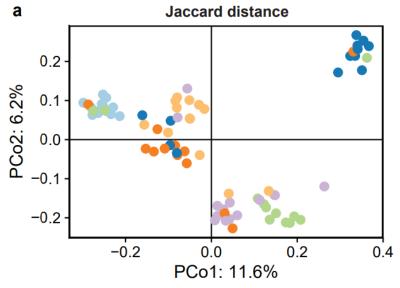
curves showing the number of unique sequence variants as a function of normalized library size

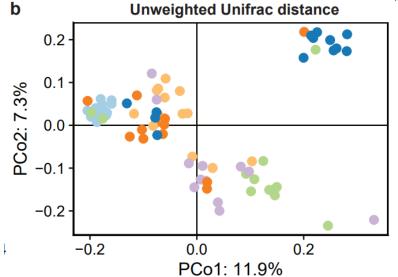
Rarefraction curve

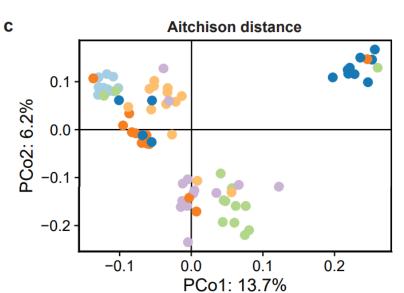


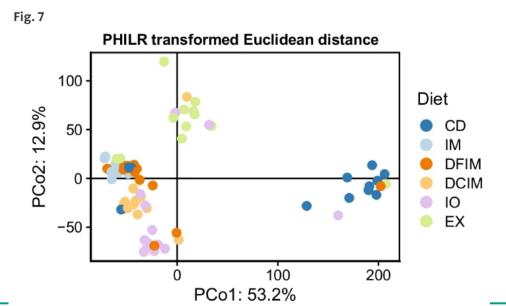




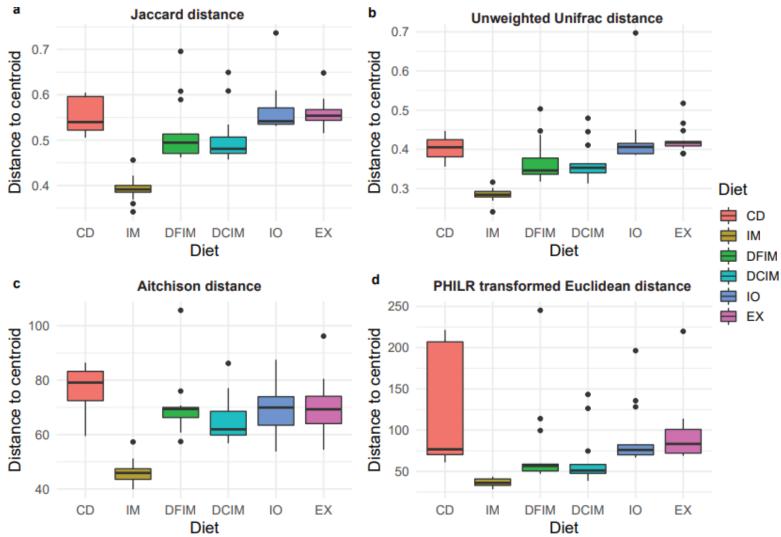












# Norwegian University of Life Sciences

68

### Tidy taxa table of top 15

```
taxa tab <- as.data.frame(taxa tab) %>%
 rownames to column("Taxa") %>%
 separate(
 Taxa,
 c("Kingdom", "Phylum", "Class", "Order", "Family", "Genus")) %>%
 Phylum = ifelse(
 is.na(Phylum)|Phylum == "NA"|grepl("uncultured|Ambiguous|metagenome", Phylum),
 Phylum),
 Class = ifelse(
 is.na(Class) Class == "NA" | grepl("uncultured | Ambiguous | metagenome", Class),
 Class),
 Order = ifelse(
 is.na(Order) | Order == "NA" | grepl("uncultured | Ambiguous | metagenome", Order),
 Class,
 Order),
 Family = ifelse(
 is.na(Family)|Family == "NA"|grepl("uncultured|Ambiguous|metagenome", Family),
 Order.
 Family),
 Genus = ifelse(
 is.na(Genus) | Genus == "NA" | grepl("uncultured | Ambiguous | metagenome", Genus),
 Family,
 Genus)) %>%
 select(-Kingdom, -(Class:Family))
taxa tab1 <- taxa tab %>%
 mutate(
 Phylum = gsub("p__", "", Phylum),
 Phylum = factor(Phylum, levels = rev(unique(Phylum))),
 Genus = gsub("g__", "", Genus),
```

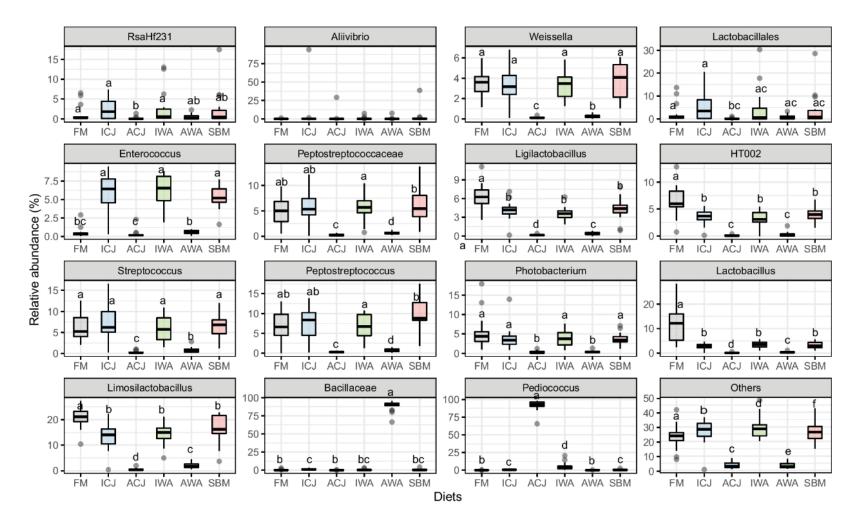


### Boxplot of individual groups

```
Boxplot to abundance of individual dominant taxa in the DIGESTA samples grouped by diets
 `{r, fig.width=10}
col <- c("grey", brewer.pal(n = 10, name = "Paired"))</pre>
mtd$Diet <- factor(mtd$Diet, levels = c("D1", "D2", "D3", "D4"))</pre>
taxa tab1$Genus2 <- reorder(taxa tab1$Genus, taxa tab1$Abundance)
taxa_boxplot_digesta <- filter(taxa_tab1, Sample_type == "digesta") %>%
 ggplot(aes(x = factor(Diet, c("D1", "D2", "D3", "D4")), y = Abundance, fill = Genus)) +
 geom boxplot(aes(fill = factor(Diet, c("D1", "D2", "D3", "D4"))), alpha = 0.5, width = 0.5) +
 labs(x = "Diets", y = "Relative abundance (%)") +
 scale fill manual(values = col) +
 facet wrap(
 \sim Genus2, nrow = 4,
 scale = "free"
 theme bw() +
 theme(##axis.text.x = element blank(),
 legend.position = "none")
ggsave("boxplot abundance digesta.tiff", width = 10, height = 6,
 units = "in", dpi = 300, compression = "lzw")
ggsave("boxplot abundance digesta.pdf", width = 10, height = 6,
 units = "in", dpi = 300)
```



### Boxplot of individual groups



### Filtering groups and perform statistic groups

```
##Flitering the dataframe into each individual taxonomic
taxa tab2 <- filter(taxa tab1, Sample type == "digesta")
taxa h <- taxa tab2 %>%
 filter(taxa tab2$Genus == "HT002")
##Filter taxa_tab2 dataframe for only the Corynebacterium
taxa_c <- taxa_tab2 %>%
 filter(taxa tab2$Genus == "Corynebacterium")
taxa leu <- taxa tab2 %>%
 filter(taxa tab2$Genus == "Leuconostoc")
taxa fur <- taxa tab2 %>%
 filter(taxa tab2$Genus == "Furfurilactobacillus")
taxa sec <- taxa tab2 %>%
 filter(taxa tab2$Genus == "Secundilactobacillus")
taxa lac1 <- taxa tab2 %>%
 filter(taxa tab2$Genus == "Lactiplantibacillus")
taxa w <- taxa tab2 %>%
 filter(taxa tab2$Genus == "Weissella")
```

```
Norwegian University of Life Sciences
```

```
Test the diet effect on individual taxonomy - Kruska-walis tes
kruskal.test(Abundance ~ Diet, data = taxa h)
#Statistics for Corynebacterium
kruskal.test(Abundance ~ Diet, data = taxa c)
#Statistics for Leuconostoc
kruskal.test(Abundance ~ Diet, data = taxa leu)
#Statistics for Furfurilactobacillus
kruskal.test(Abundance ~ Diet, data = taxa fur)
#Statistics for <u>Secundilactobacillus</u>
kruskal.test(Abundance ~ Diet, data = taxa sec)
#Statistics for Lactiplantibacillus
kruskal.test(Abundance ~ Diet, data = taxa lac1)
#Statistics for Weissella
```

### 7. LEfSe



Segata et al. Genome Biology 2011, 12:R60 http://genomebiology.com/2011/11/6/R60



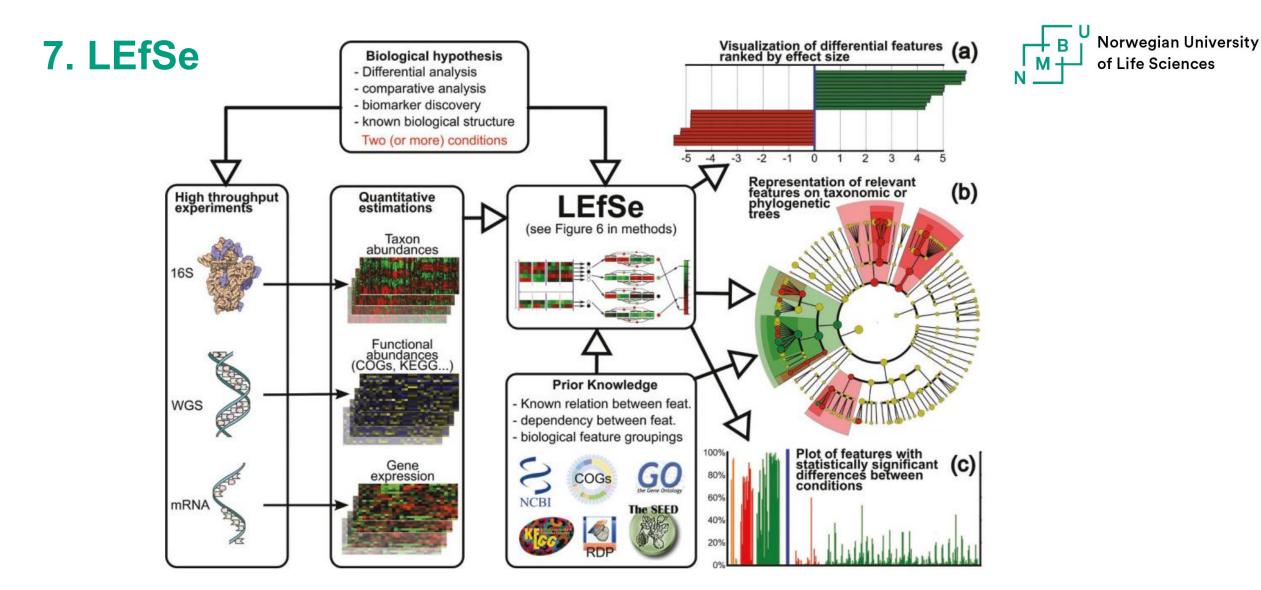
METHOD Open Access

# Metagenomic biomarker discovery and explanation

Nicola Segata<sup>1</sup>, Jacques Izard<sup>2,3</sup>, Levi Waldron<sup>1</sup>, Dirk Gevers<sup>4</sup>, Larisa Miropolsky<sup>1</sup>, Wendy S Garrett<sup>5,6,7</sup> and Curtis Huttenhower<sup>1\*</sup>

#### Abstract

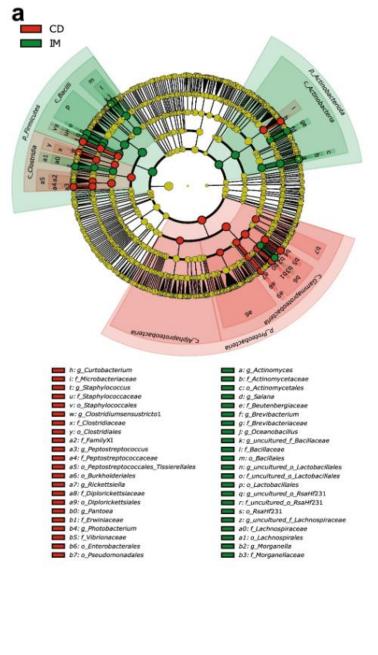
This study describes and validates a new method for metagenomic biomarker discovery by way of class comparison, tests of biological consistency and effect size estimation. This addresses the challenge of finding organisms, genes, or pathways that consistently explain the differences between two or more microbial communities, which is a central problem to the study of metagenomics. We extensively validate our method on several microbiomes and a convenient online interface for the method is provided at http://huttenhower.sph. harvard.edu/lefse/.

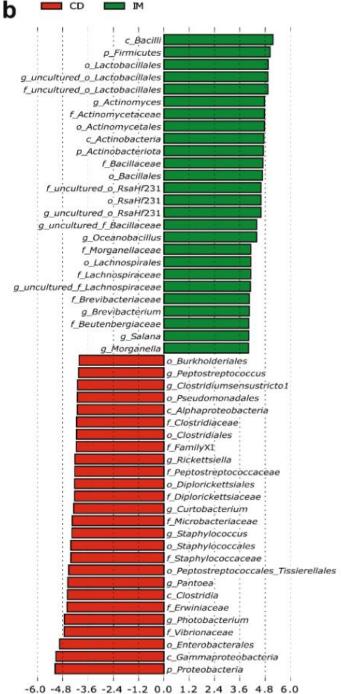


http://galaxy.biobakery.org/

FOODS PNORWAY 73

### 7. LEfSe





LDA SCORE (log 10)



Weththasinghe, P., et al. (2022)

### 8. Metabolic pathways



