

# Dada2 pipeline

by BiomCare

10/12/2019

Customer	SEGES
Customer ID	DA00201-19
Project ID	Differences between bacetria and fungi in plowed and unplowed clay and sand samples
Sample Type	Soil

In this Data Processing Report, you will find information regarding DNA sequencing, data quality evaluation and filtering, and microbiome profiling.

All your data, illustrations, supplementary files and reports are now available for download in your private project folder on BiomCare's cloud. Please refer to the supplementary document How to navigate your BiomCare folder for details on how to find specific files and for instructions on how to interpret illustrations that are not introduced in this report.

## Summary

Paired-end 16S rRNA gene amplicon metagenomic sequencing was performed, targeting 10.000 reads per sample (details on sequencing are available below). The performed amplicon sequencing resulted in a mean read count of 17758 reads across samples.

Evaluation of the quality of the raw data show a base quality score (phred score) above [28] across read lengths. No sample (single fastq file) has reads flagged as poor quality. The GC content is on average 56.54167 (min: 56 , max: 57). Reads across samples has a mean read length of 301bp.

The raw data does not contain adaptors/primers due to the use of customized primers for library preparation, and therefore no removal of non-biological sequences was performed. Quality filtering was performed using DADA2 and removed on average 844.6667 (min: 255, max: 2949) reads per sample, while denoising of reads removed on average 2369.292 reads across samples. The read count following filtering and microbiome profiling was above the target threshold of 10.000 reads per sample for \*\*\* samples. The remaining samples (n= ) had a low read-count following processing and we recommend removing these from further analysis. Inspection of the included quality control samples, as well as the mock community, indicated no contamination and a successful sequencing, filtering and microbiome profiling.

## Table of contents

- Summary
- Table of contents
- Data generation using Massive Parallel Sequencing

- Evaluation of raw data quality
- Quality filtering
  - Removal of adapter/primer sequences
  - Read quality and length filtering.
  - Calculate error rates
  - Sample Inference and Denoising
  - Remove chimeras
  - Read counts
- Software, settings, thresholds and reference data used

## Data generation using Massive Parallel Sequencing

Sequencing is performed using polymerase chain reaction (PCR) for DNA amplification of the selected variable region(s) of the bacterial 16S rRNA gene. The product from PCR is normalized using the SequelPrep Normalization Kit, followed by sequencing on an Illumina MiSeq platform using v3 chemistry and multiplexing of up to 384 samples per run. This results in 2x300bp paired-end reads, that are demultiplexed while allowing for no mismatches in the index sequences.

## Evaluation of raw data quality

BiomCare use FastQC and DADA2 to evaluate the quality of the raw sequencing data. BiomCare use the results from FastQC and DADA2 to evaluate the quality of the sequencing data, and to design appropriate quality filtering settings and steps. If results from the quality evaluation indicate an issue with data generation, we bring this back to our sequencing facility (if BiomCare has generated the data) or to you (if you have provided us with raw sequencing data). The summary section above contains the key information from the data quality evaluation.

If you wish to further evaluate the results of the data quality assessment, please refer to the supportive documents in your BiomCare folder (location: 2\_Quality\_filtering\_and\_profiling/Illustrations) and the guide on how to read the illustrations found in the document Report Attachment 2.

Below you find a description of BiomCares data quality filtering. If data contain adaptors, adaptor removal is performed as part of the quality filtering.

## Quality filtering

Quality filtering steps

- Removal of adapter/primer sequences
- Read quality and length filtering

## Removal of adapter/primer sequences

Our sequencing facility uses custom primers wherever possible, leading to no sequence trimming of adapter/primer data needed for those samples. However, for those other samples, we always evaluate what measures are required to remove adaptor sequences from their next generation sequencing data. For some sequencing platforms, this is done during demultiplexing, however for

others it is necessary to remove the adaptor sequences during quality filtering. It is good practices to evaluate the presence of adaptor sequences and remove any remaining adaptor sequences before further processing. This step is performed using the software Cutadapt.

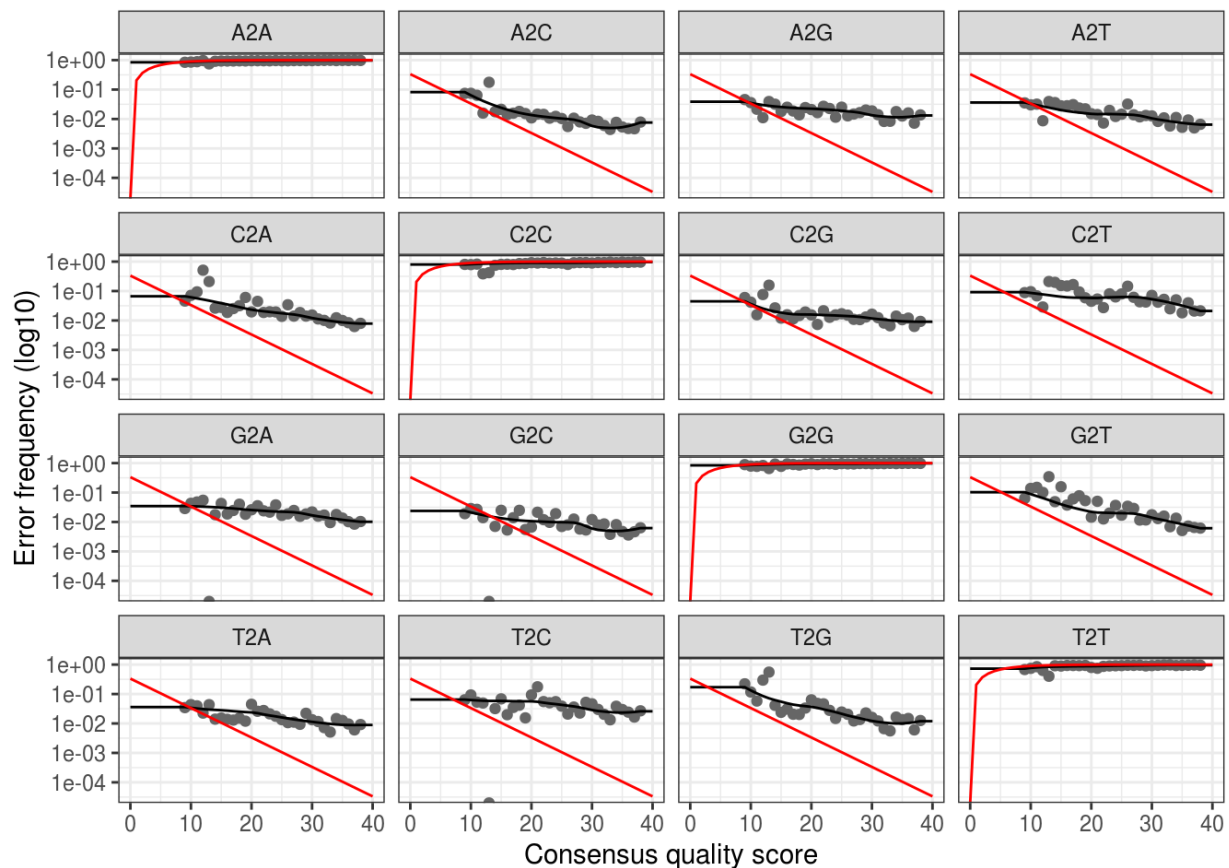
## Read quality and length filtering

When bases are called based on the data obtained from the sequencing platform, each base is annotated with a quality score. The quality score of the called bases (often in the form of phred scores) is accessed and used to remove both low-quality reads, and low-quality bases at the ends of reads. Reads that are shorter than a defined length threshold following quality filtering are removed. This step is performed using DADA2.

We do this for a number of reasons: a) to remove noise that may limited the information obtainable from the quality evaluation, b) to see if quality issues are easily handled with standard processing and c) to evaluate first-choice quality settings.

## Calculate error rates

The errors made in the sequencing process generates a specific profile for each dataset. In order for to denoise the sequences, an error model needs to be built. This is done using machine learning to estimate the error rates for each possible nucleotide transition for this specific dataset. This step is performed using DADA2 and the error rates can be seen below in **Figure 1**.



**Figure 1: Error rates for all nucleotide transitions.**

## Sample Inference and denoising

DADA2 uses the error rates calculated in the previous step in order to correct for any sequencing errors present. This denoising step allows for sequence inference, establishing what reads belong to the same sequence. DADA2 performs these steps as one step in its signature Divisive Amplicon Denoising Algorithm, leaving only those sequences which are likely to be truly present in the samples.

## Remove chimeras

DNA sequencing data contain a low number of sequences called chimeras. Chimera sequences are artifacts that are formed when two or more biological sequences are incorrectly joined. Because these sequences do not represent a true microorganism, they must be removed before further data processing. Chimeras are identified by first locating the two parent sequences of the chimeras using the Needleman-Wunsch global algorithm. Finding two parent sequences which overlap exactly the child sequence from the left and right indicates the finding of a chimera. The identified chimeras are then removed.

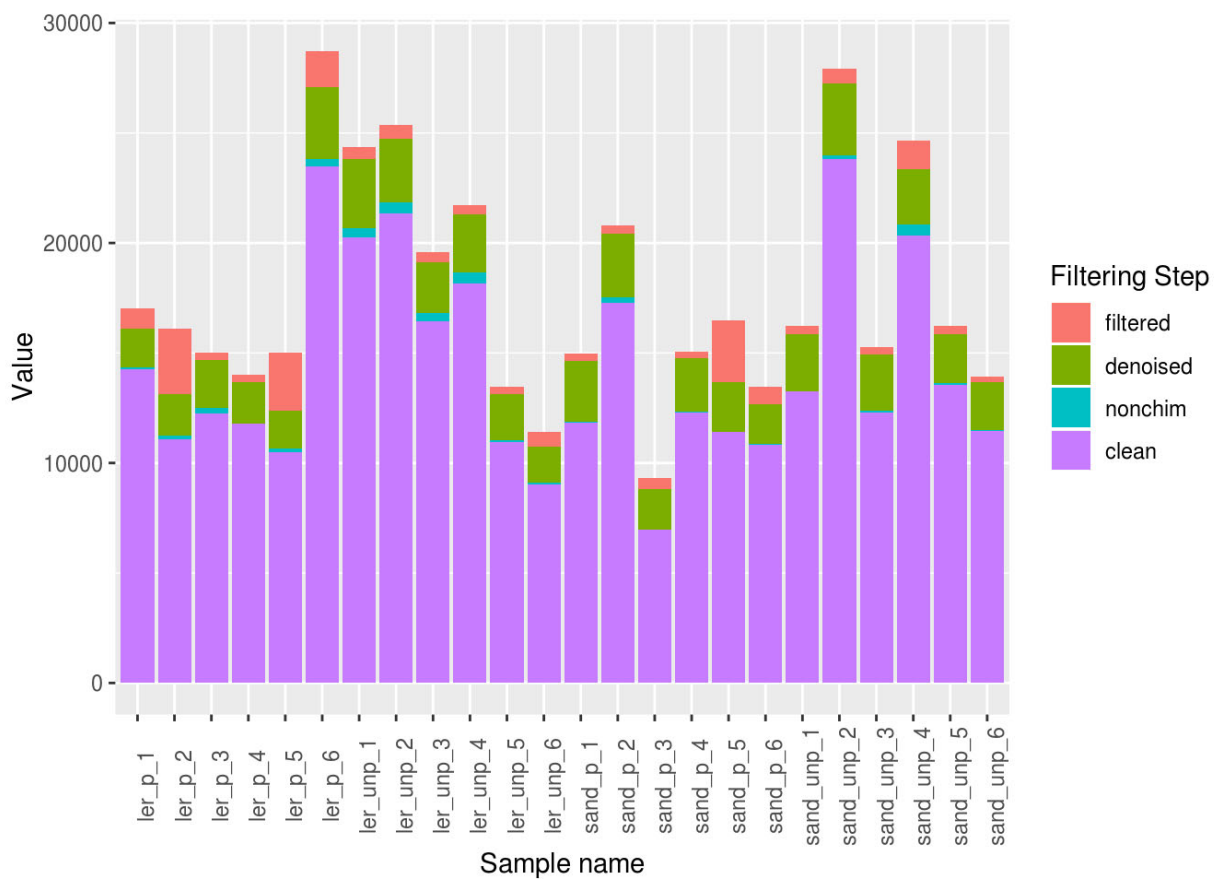
## Read counts

For each sample, reads may be removed at each of the quality filtering steps. Inspecting the number of reads filtered provide important information on the data quality. Table 1 below gives an easy overview of the data at each step through quality filtering. The table include an overview of the reads across samples in the raw and cleaned data (first and last row of the table).

**Table 1: Summary statistics for the data throughout filtering.** The first and last row provide summary information on read counts in the raw sequencing data and the cleaned data following the full microbiome profiling process (merged reads), respectively. The three rows in-between (filtered, denoised and merged) are the summary information for the number of reads removed at each of these steps across all samples.

	Min.	1st Q.	Median	Mean	3rd Q.	Max.
<b>Input</b>	9332	14737.00	16164.0	17758.00	21026.00	28724
<b>Filtered</b>	255	349.20	498.0	844.70	817.00	2949
<b>Denoised</b>	1637	1914.00	2286.0	2369.00	2676.00	3261
<b>Nonchim</b>	0	43.75	116.5	178.12	284.25	487
<b>Clean</b>	6953	11310.00	12286.0	14366.00	17491.00	23809

A supportive file is provided for visual inspection of read counts in each sample at each filtering step (see file: Read\_counts\_per\_sample.pdf). Please see Figure 2 below for an introduction to the barplots found in the supportive pdf.



**Figure 2: Summary of reads across samples in the project.** Each bar show reads for each sample. The “cleaned” reads are the number of reads (merged if paired-end data) that remain after processing the raw data and calling the microbiome profiles. Filtered is the number of reads removed during read filtering (removed due to low read quality or short length). Denoised is the number of reads removed due to duplication.

After data processing, some samples are likely to have a very low read count. This will in most instances be caused by low DNA concentration or quality of the sample that was sequenced. To support an effective evaluation of the samples with a low read count, stacked barplots for these samples are shown in a separate PDF file (file\_name). Samples were selected based on a read count below 5000 after data processing, as this is the number of reads that will be available for further analysis.

## Software, settings, thresholds and reference data used

**Table 2: Reference and parameter values**

Parameters	Values
Truncation quality	5
Truncation length	F:250, R:200
Left trimming length	F:5, R:5

Parameters	Values
Right trimming length	0
Max. length	Inf
Min. length	20
Max. #N	0
Min. Q	0
Max. Expected Errors	F:2, R:2
Remove PhiX	TRUE
Low complexity filter	0
Reference Data	RDP training set 16

**Table 3: Software and package versions**

	version
<b>OS</b>	Ubuntu 16.04.6 LTS
<b>R</b>	3.6.1
<b>colorspace</b>	1.4-1
<b>hwriter</b>	1.3.2
<b>ellipsis</b>	0.3.0
<b>rprojroot</b>	1.3-2
<b>XVector</b>	0.26.0
<b>GenomicRanges</b>	1.38.0
<b>fs</b>	1.3.1
<b>rstudioapi</b>	0.10
<b>farver</b>	2.0.1
<b>remotes</b>	2.1.0
<b>lubridate</b>	1.7.4
<b>xml2</b>	1.2.2

	version
<b>codetools</b>	0.2-16
<b>splines</b>	3.6.1
<b>knitr</b>	1.26
<b>pkgload</b>	1.0.2
<b>polyclip</b>	1.10-0
<b>zeallot</b>	0.1.0
<b>ade4</b>	1.7-13
<b>Rsamtools</b>	2.2.1
<b>broom</b>	0.5.2
<b>cluster</b>	2.1.0
<b>dbplyr</b>	1.4.2
<b>compiler</b>	3.6.1
<b>httr</b>	1.4.1
<b>backports</b>	1.1.5
<b>assertthat</b>	0.2.1
<b>Matrix</b>	1.2-17
<b>lazyeval</b>	0.2.2
<b>cli</b>	1.1.0
<b>tweenr</b>	1.0.1
<b>prettyunits</b>	1.0.2
<b>htmltools</b>	0.4.0
<b>tools</b>	3.6.1
<b>igraph</b>	1.2.4.2
<b>gtable</b>	0.3.0
<b>glue</b>	1.3.1

	version
<b>GenomeInfoDbData</b>	1.2.2
<b>reshape2</b>	1.4.3
<b>ShortRead</b>	1.44.0
<b>Biobase</b>	2.46.0
<b>cellranger</b>	1.1.0
<b>vctrs</b>	0.2.0
<b>Biostrings</b>	2.54.0
<b>multtest</b>	2.42.0
<b>ape</b>	5.3
<b>nlme</b>	3.1-140
<b>iterators</b>	1.0.12
<b>xfun</b>	0.11
<b>ps</b>	1.3.0
<b>testthat</b>	2.3.1
<b>rvest</b>	0.3.5
<b>lifecycle</b>	0.1.0
<b>devtools</b>	2.2.1
<b>zlibbioc</b>	1.32.0
<b>MASS</b>	7.3-51.4
<b>scales</b>	1.1.0
<b>hms</b>	0.5.2
<b>parallel</b>	3.6.1
<b>SummarizedExperiment</b>	1.16.0
<b>biomformat</b>	1.14.0
<b>rhdf5</b>	2.30.1



	version
<b>RColorBrewer</b>	1.1-2
<b>curl</b>	4.3
<b>yaml</b>	2.2.0
<b>memoise</b>	1.1.0
<b>latticeExtra</b>	0.6-28
<b>stringi</b>	1.4.3
<b>highr</b>	0.8
<b>desc</b>	1.2.0
<b>S4Vectors</b>	0.24.1
<b>foreach</b>	1.4.7
<b>permute</b>	0.9-5
<b>BiocGenerics</b>	0.32.0
<b>pkgbuild</b>	1.0.6
<b>zip</b>	2.0.4
<b>BiocParallel</b>	1.20.0
<b>GenomeInfoDb</b>	1.22.0
<b>rlang</b>	0.4.2
<b>pkgconfig</b>	2.0.3
<b>matrixStats</b>	0.55.0
<b>bitops</b>	1.0-6
<b>evaluate</b>	0.14
<b>lattice</b>	p.20-38
<b>Rhdf5lib</b>	1.8.0
<b>labeling</b>	0.3
<b>GenomicAlignments</b>	1.22.1

	version
<b>processx</b>	3.4.1
<b>tidyselect</b>	0.2.5
<b>plyr</b>	1.8.4
<b>magrittr</b>	1.5
<b>R6</b>	2.4.1
<b>IRanges</b>	2.20.1
<b>generics</b>	0.0.2
<b>DelayedArray</b>	0.12.0
<b>DBI</b>	1.0.0
<b>pillar</b>	1.4.2
<b>haven</b>	2.2.0
<b>withr</b>	2.1.2
<b>mgcv</b>	1.8-28
<b>survival</b>	2.44-1.1
<b>RCurl</b>	1.95-4.12
<b>modelr</b>	0.1.5
<b>crayon</b>	1.3.4
<b>rmarkdown</b>	1.18
<b>usethis</b>	1.5.1
<b>readxl</b>	1.3.1
<b>callr</b>	3.3.2
<b>vegan</b>	2.5-6
<b>webshot</b>	0.5.2
<b>reprex</b>	0.3.0
<b>digest</b>	0.6.23

	version
<b>RcppParallel</b>	4.4.4
<b>stats4</b>	3.6.1
<b>munsell</b>	0.5.0
<b>viridisLite</b>	0.3.0
<b>sessioninfo</b>	1.1.1