

# DADA2 analysis of 16S rRNA gene amplicon sequencing reads

Jianshu Zhao

2021-04-30

```
knitr::opts_chunk$set(echo = TRUE)
knitr::opts_chunk$set(fig.width = 10)
knitr::opts_chunk$set(tidy.opts=list(width.cutoff=60),tidy=TRUE)
```

## Introduction

Implementing DADA2 pipeline for resolving sequence variants from 16S rRNA gene amplicon *paired-end* sequencing reads, adopting the tutorial from <https://benjjneb.github.io/dada2/tutorial.html> and [https://benjjneb.github.io/dada2/bigdata\\_paired.html](https://benjjneb.github.io/dada2/bigdata_paired.html) with minor adjustments. This report captures all the workflow steps necessary to reproduce the analysis.

## Load R packages:

```
##
## Attaching package: 'RcppParallel'

## The following object is masked from 'package:Rcpp':
##
##      LdFlags
## [1] '1.18.0'
## 24
## $BAND_SIZE
## [1] 16
##
## $DETECT_SINGLETONS
## [1] FALSE
##
## $GAP_PENALTY
## [1] -8
##
## $GAPLESS
## [1] TRUE
##
## $GREEDY
## [1] TRUE
##
## $HOMOPOLYMER_GAP_PENALTY
## NULL
##
## $KDIST_CUTOFF
## [1] 0.42
```

```

##
## $MATCH
## [1] 5
##
## $MAX_CLUST
## [1] 0
##
## $MAX_CONSIST
## [1] 10
##
## $MIN_ABUNDANCE
## [1] 1
##
## $MIN_FOLD
## [1] 1
##
## $MIN_HAMMING
## [1] 1
##
## $MISMATCH
## [1] -4
##
## $OMEGA_A
## [1] 1e-40
##
## $OMEGA_C
## [1] 1e-40
##
## $OMEGA_P
## [1] 1e-04
##
## $PSEUDO_ABUNDANCE
## [1] Inf
##
## $PSEUDO_PREVALENCE
## [1] 2
##
## $SSE
## [1] 2
##
## $USE_KMERS
## [1] TRUE
##
## $USE_QUALS
## [1] TRUE
##
## $VECTORIZED_ALIGNMENT
## [1] TRUE

## /storage/coda1/p-ktk3/0/jzhao399/rich_project_bio-konstantinidis/apps/dada2_wrap
## per/scripts

```

Get list of input fastq files.

```
# Variable 'input.path' containing path to input fastq files
# directory is inherited from wrapper script dada2_cli.r.

input.file.list <- grep("*fastq", list.files(input.path), value = T)
# input.path <- normalizePath('input/')

# List of input files

# Sort ensures forward/reverse reads are in same order
fnFs <- sort(grep("_R1.*\\.fastq", list.files(input.path), value = T))
fnRs <- sort(grep("_R2.*\\.fastq", list.files(input.path), value = T))

# Extract sample names, allowing variable filenames; e.g.
# *_R1[_001].fastq.gz]
sample.names <- gsub("_R1.*\\.fastq(\\.gz)?", "", fnFs, perl = T)
sample.namesR <- gsub("_R2.*\\.fastq(\\.gz)?", "", fnRs, perl = T)
if (!identical(sample.names, sample.namesR)) stop("Forward and reverse files do not match.")

# Specify the full path to the fnFs and fnRs
fnFs <- file.path(input.path, fnFs)
fnRs <- file.path(input.path, fnRs)
```

Generate quality plots for FWD and REV reads and store in Read\_QC folder.

```
# Create output folder NOTE: variable 'output.dir' containing
# name of output folder is inherited from wrapper script
# dada2_cli.r.
cwd <- getwd()

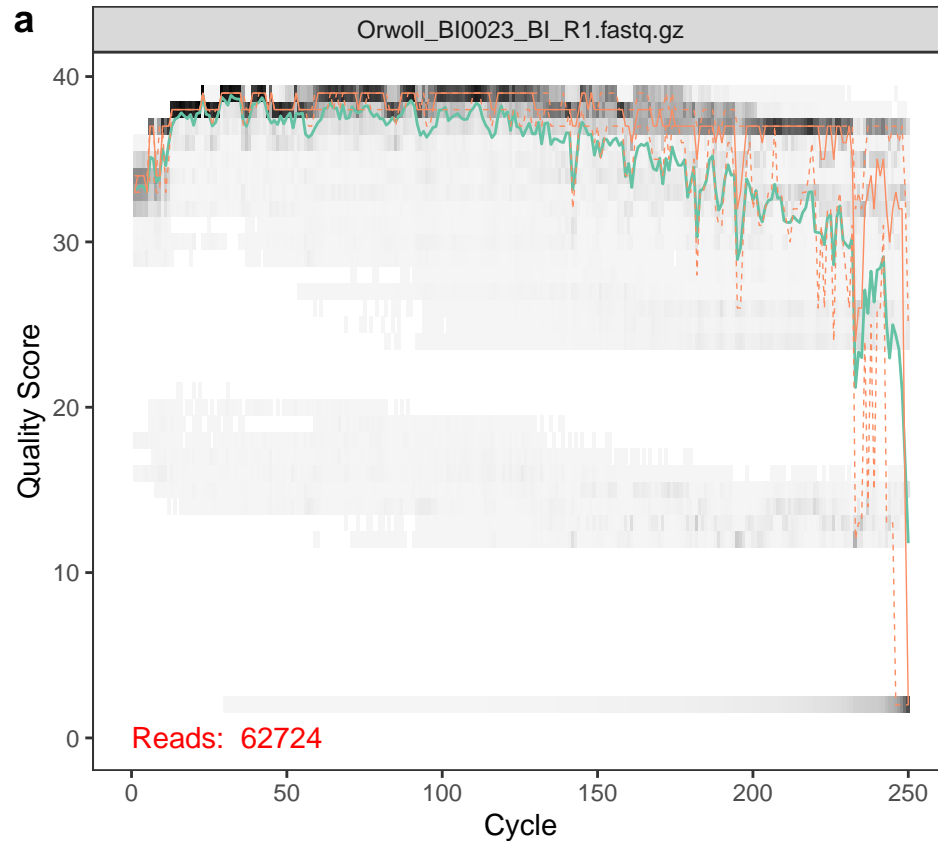
readQC.folder <- file.path(cwd, output.dir, "Read_QC")
ifelse(!dir.exists(readQC.folder), dir.create(readQC.folder,
  recursive = TRUE), FALSE)

## [1] TRUE

# Generate plots and save to folder in multi-page pdf

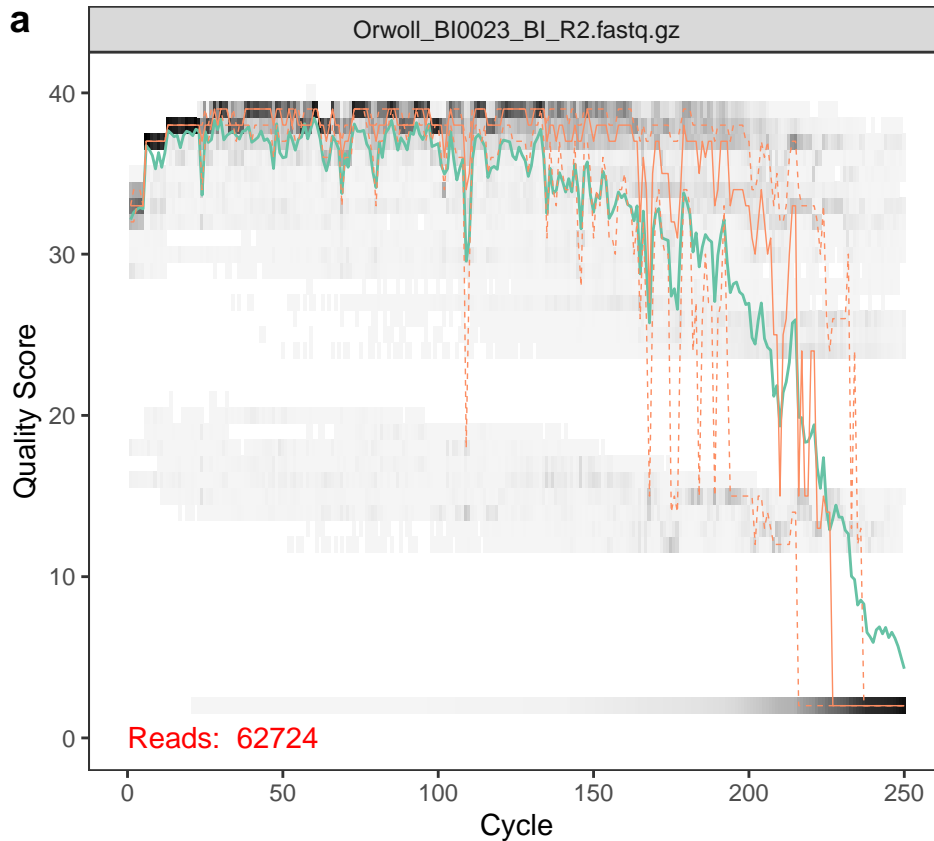
# Forward reads
fwd.qc.plots.list <- list()
for (i in 1:length(fnFs)) {
  fwd.qc.plots.list[[i]] <- plotQualityProfile(fnFs[i])
  rm(i)
}

## plot the quality plot of forward reads of first sample
## fwd.qc.plots.list[[1]] Save to file
mF <- marrangeGrob(fwd.qc.plots.list, ncol = 2, nrow = length(fnFs)/2 +
  1, top = NULL)
ggsave(paste0(readQC.folder, "/FWD_read_plot.pdf"), mF, width = 7,
  height = 2.5 * (length(fnFs)/2 + 1), device = "pdf", limitsize = FALSE)
plot_grid(fwd.qc.plots.list[[1]], ncol = 2, labels = "auto")
```



```
rm(fwd.qc.plots.list)

# Reverse reads
rev.qc.plots.list <- list()
for (i in 1:length(fnRs)) {
  rev.qc.plots.list[[i]] <- plotQualityProfile(fnRs[i])
  rm(i)
}
## plot the quality plot of forward reads of first sample
## rev.qc.plots.list[[1]]
mR <- marrangeGrob(rev.qc.plots.list, ncol = 2, nrow = length(fnRs)/2 +
  1, top = NULL)
ggsave(paste0(readQC.folder, "/REV_read_plot.pdf"), mR, width = 7,
  height = 2.5 * (length(fnRs)/2 + 1), device = "pdf", limitsize = FALSE)
plot_grid(rev.qc.plots.list[[1]], ncol = 2, labels = "auto")
```



```
rm(rev.qc.plots.list)
```

## Trim and filter reads.

```
# Create filtered_input/ subdirectory for storing filtered
# fastq reads
filt_path <- file.path(cwd, output.dir, "filtered_input")
ifelse(!dir.exists(filt_path), dir.create(filt_path, recursive = TRUE),
      FALSE)

## [1] TRUE

# Define filenames for filtered input files
filtFs <- file.path(filt_path, paste0(sample.names, "_F_filt.fastq.gz"))
filtRs <- file.path(filt_path, paste0(sample.names, "_R_filt.fastq.gz"))

fnFs.noprimer <- file.path(filt_path, paste0(sample.names, "noprimer_F_filt.fastq.gz"))
fnRs.noprimer <- file.path(filt_path, paste0(sample.names, "noprimer_R_filt.fastq.gz"))

# remove primers
if (!is.null(Fwd_pr) & !is.null(Rvs_pr)) {
  removePrimers(fnFs, fnFs.noprimer, primer.fwd = Fwd_pr, trim.fwd = TRUE,
    orient = TRUE, verbose = TRUE)
  removePrimers(fnRs, fnRs.noprimer, primer.fwd = Rvs_pr, trim.fwd = TRUE,
    orient = TRUE, verbose = TRUE)
} else {
  fnFs.noprimer <- fnFs
```

```

fnRs.noprimer <- fnRs
}

# Filter the forward and reverse reads: Note that: 1. Reads
# are both truncated and then filtered using the maxEE
# expected errors algorithm from UPARSE. 2. Reverse reads
# are truncated to shorter lengths than forward since they
# are much lower quality. 3. _Both_ reads must pass for the
# read pair to be output. 4. Output files are compressed by
# default.

rd.counts <- as.data.frame(filterAndTrim(fnFs.noprimer, filtFs,
  fnRs.noprimer, filtRs, truncLen = c(200, 190), maxN = 0,
  maxEE = c(1, 2), truncQ = 10, rm.phix = TRUE, compress = TRUE,
  multithread = threads, matchIDs = TRUE))
# Table of before/after read counts
rd.counts$ratio <- round(rd.counts$reads.out/rd.counts$reads.in,
  digits = 2)
rd.counts

##
## reads.in reads.out ratio
## Orwoll_BI0023_BI_R1.fastq.gz 62724 53296 0.85
## Orwoll_BI0056_BI_R1.fastq.gz 55342 47545 0.86
## Orwoll_BI0131_BI_R1.fastq.gz 55144 46797 0.85
## Orwoll_BI0153_BI_R1.fastq.gz 44610 39510 0.89
## Orwoll_BI0215_BI_R1.fastq.gz 48227 40256 0.83
## Orwoll_BI0353_BI_R1.fastq.gz 54271 49323 0.91

# Write rd.counts table to file in readQC.folder
write.table(rd.counts, paste0(readQC.folder, "/Read_counts_after_filtering.tsv"),
  sep = "\t", quote = F, eol = "\n", col.names = NA)

```

## Learn the error rates.

The DADA2 algorithm depends on a parametric error model (err) and every amplicon dataset has a different set of error rates. The learnErrors method learns the error model from the data, by alternating estimation of the error rates and inference of sample composition until they converge on a jointly consistent solution. As in many optimization problems, the algorithm must begin with an initial guess, for which the maximum possible error rates in this data are used (the error rates if only the most abundant sequence is correct and all the rest are errors).

```

set.seed(100)
# Filtered forward reads
errF <- learnErrors(filtFs, nread = 1e+06, multithread = threads,
  randomize = TRUE)

## Warning in learnErrors(filtFs, nread = 1e+06, multithread = threads, randomize
## = TRUE): The nreads parameter is DEPRECATED. Please update your code with the
## nbases parameter.

## 55345400 total bases in 276727 reads from 6 samples will be used for learning the error rates.

# Filtered reverse reads
errR <- learnErrors(filtRs, nread = 1e+06, multithread = threads,
  randomize = TRUE)

## Warning in learnErrors(filtRs, nread = 1e+06, multithread = threads, randomize

```

```
## = TRUE): The nreads parameter is DEPRECATED. Please update your code with the
## nbases parameter.

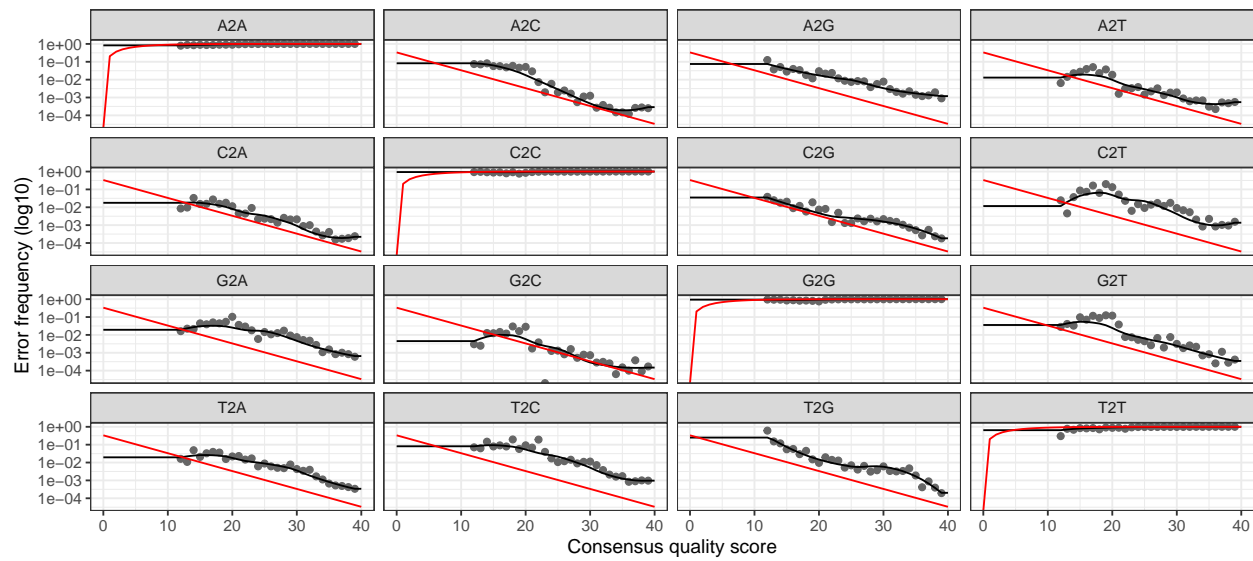
## 52578130 total bases in 276727 reads from 6 samples will be used for learning the error rates.
```

```
# Visualize the estimated error rates Forward
```

```
plotErrors(errF, nominalQ = TRUE)
```

```
## Warning: Transformation introduced infinite values in continuous y-axis
```

```
## Warning: Transformation introduced infinite values in continuous y-axis
```



```
# Save to file
```

```
ggsave(paste0(readQC.folder, "/Error_rates_per_sample_FWD.pdf"),
  plotErrors(errF, nominalQ = TRUE), device = "pdf")
```

```
## Saving 10 x 4.5 in image
```

```
## Warning: Transformation introduced infinite values in continuous y-axis
```

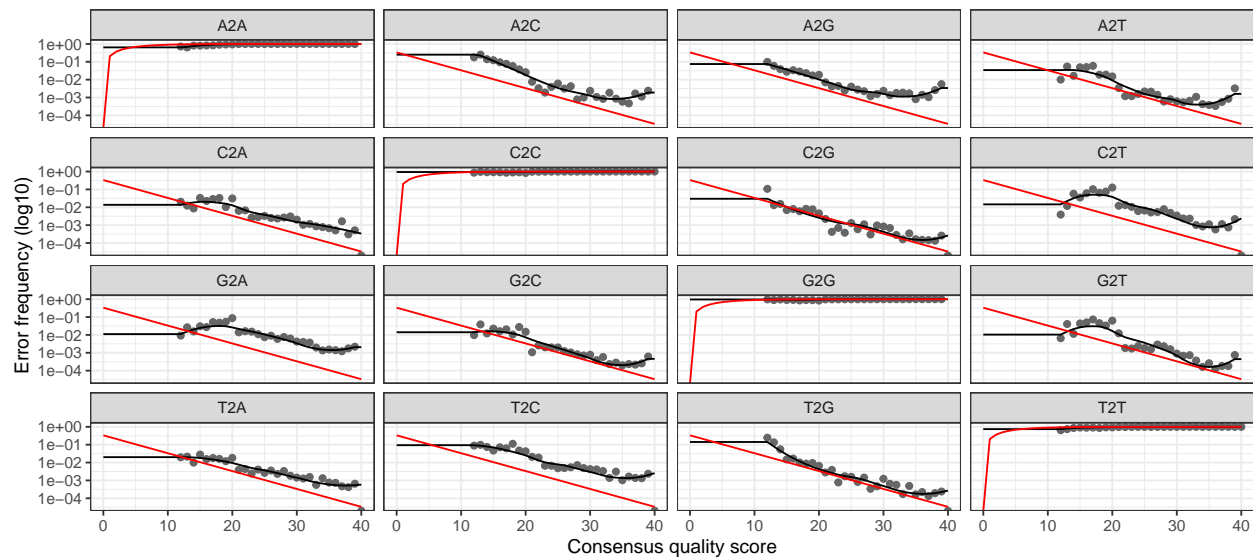
```
## Warning: Transformation introduced infinite values in continuous y-axis
```

```
# Reverse
```

```
plotErrors(errR, nominalQ = TRUE)
```

```
## Warning: Transformation introduced infinite values in continuous y-axis
```

```
## Warning: Transformation introduced infinite values in continuous y-axis
```



```
# Save to file
ggsave(paste0(readQC.folder, "/Error_rates_per_sample_REV.pdf"),
       plotErrors(errR, nominalQ = TRUE), device = "pdf")
```

```
## Saving 10 x 4.5 in image
```

```
## Warning: Transformation introduced infinite values in continuous y-axis
```

```
## Warning: Transformation introduced infinite values in continuous y-axis
```

The error rates for each possible transition (eg. A->C, A->G, ...) are shown. Points are the observed error rates for each consensus quality score. The black line shows the estimated error rates after convergence. The red line shows the error rates expected under the nominal definition of the Q-value. If the black line (the estimated rates) fits the observed rates well, and the error rates drop with increased quality as expected, then everything looks reasonable and can proceed with confidence.

## Infer Sequence Variants

This step consists of dereplication, sample inference, and merging of paired reads

Dereplication combines all identical sequencing reads into into “unique sequences” with a corresponding “abundance”: the number of reads with that unique sequence. DADA2 retains a summary of the quality information associated with each unique sequence. The consensus quality profile of a unique sequence is the average of the positional qualities from the dereplicated reads. These quality profiles inform the error model of the subsequent denoising step, significantly increasing DADA2’s accuracy.

The sample inference step performs the core sequence-variant inference algorithm to the dereplicated data.

Spurious sequence variants are further reduced by merging overlapping reads. The core function here is `mergePairs`, which depends on the forward and reverse re.samples being in matching order at the time they were dereplicated.

```
# Sample inference of dereplicated reads, and merger of
# paired-end reads
if (!pool.samples == "TRUE") {
  mergers <- vector("list", length(sample.names))
  names(mergers) <- sample.names
  names(filtFs) <- sample.names
  names(filtRs) <- sample.names
```



```

for (sam in sample.names) {
  cat("Processing:", sam, "\n")
  derepF <- derepFastq(filtFs[[sam]])
  ddF <- dada(derepF, err = errF, multithread = threads)
  derepR <- derepFastq(filtRs[[sam]])
  ddR <- dada(derepR, err = errR, multithread = threads)
  merger <- mergePairs(ddF, derepF, ddR, derepR)
  mergers[[sam]] <- merger
}
cat(pool.samples)
rm(derepF)
rm(derepR)
} else {
  derepFs <- derepFastq(filtFs, verbose = TRUE)
  derepRs <- derepFastq(filtRs, verbose = TRUE)
  names(derepFs) <- sample.names
  names(derepRs) <- sample.names
  dadaFs <- dada(derepFs, err = errF, multithread = threads,
    pool = TRUE)
  dadaRs <- dada(derepRs, err = errR, multithread = threads,
    pool = TRUE)
  mergers <- mergePairs(dadaFs, derepFs, dadaRs, derepRs, verbose = TRUE)
  cat(pool.samples)
  rm(derepFs)
  rm(derepRs)
}

```

```

## Dereplicating sequence entries in Fastq file: /storage/coda1/p-ktk3/0/jzhao399/rich_project_bio-kons
## Encountered 10285 unique sequences from 53296 total sequences read.
## Dereplicating sequence entries in Fastq file: /storage/coda1/p-ktk3/0/jzhao399/rich_project_bio-kons
## Encountered 6843 unique sequences from 47545 total sequences read.
## Dereplicating sequence entries in Fastq file: /storage/coda1/p-ktk3/0/jzhao399/rich_project_bio-kons
## Encountered 7450 unique sequences from 46797 total sequences read.
## Dereplicating sequence entries in Fastq file: /storage/coda1/p-ktk3/0/jzhao399/rich_project_bio-kons
## Encountered 7331 unique sequences from 39510 total sequences read.
## Dereplicating sequence entries in Fastq file: /storage/coda1/p-ktk3/0/jzhao399/rich_project_bio-kons
## Encountered 7861 unique sequences from 40256 total sequences read.
## Dereplicating sequence entries in Fastq file: /storage/coda1/p-ktk3/0/jzhao399/rich_project_bio-kons
## Encountered 7184 unique sequences from 49323 total sequences read.
## Dereplicating sequence entries in Fastq file: /storage/coda1/p-ktk3/0/jzhao399/rich_project_bio-kons
## Encountered 15600 unique sequences from 53296 total sequences read.
## Dereplicating sequence entries in Fastq file: /storage/coda1/p-ktk3/0/jzhao399/rich_project_bio-kons
## Encountered 11149 unique sequences from 47545 total sequences read.
## Dereplicating sequence entries in Fastq file: /storage/coda1/p-ktk3/0/jzhao399/rich_project_bio-kons
## Encountered 12109 unique sequences from 46797 total sequences read.

```

```
## Dereplicating sequence entries in Fastq file: /storage/coda1/p-ktk3/0/jzhao399/rich_project_bio-kons
## Encountered 11085 unique sequences from 39510 total sequences read.
## Dereplicating sequence entries in Fastq file: /storage/coda1/p-ktk3/0/jzhao399/rich_project_bio-kons
## Encountered 11940 unique sequences from 40256 total sequences read.
## Dereplicating sequence entries in Fastq file: /storage/coda1/p-ktk3/0/jzhao399/rich_project_bio-kons
## Encountered 12168 unique sequences from 49323 total sequences read.
## 6 samples were pooled: 276727 reads in 40449 unique sequences.
## 6 samples were pooled: 276727 reads in 66548 unique sequences.
## 49613 paired-reads (in 704 unique pairings) successfully merged out of 51976 (in 1545 pairings) input
## 45841 paired-reads (in 515 unique pairings) successfully merged out of 47042 (in 1038 pairings) input
## 44226 paired-reads (in 578 unique pairings) successfully merged out of 46115 (in 1279 pairings) input
## 37388 paired-reads (in 667 unique pairings) successfully merged out of 38980 (in 1328 pairings) input
## 37812 paired-reads (in 643 unique pairings) successfully merged out of 39519 (in 1367 pairings) input
## 46633 paired-reads (in 463 unique pairings) successfully merged out of 48652 (in 1114 pairings) input
## TRUE
```

## Construct sequence table

We can now construct a “sequence table” of our samples, a higher-resolution version of the “OTU table” produced by classical methods:

```
seqtab <- makeSequenceTable(mergers)
dim(seqtab)
```

```
## [1]    6 1452
```

```
# Inspect distribution of sequence lengths
table(nchar(getSequences(seqtab)))
```

```
##
## 252 253 254
## 87 1357 8
```

```
# The sequence table is a matrix with rows corresponding to
# (and named by) the samples, and columns corresponding to
# (and named by) the sequence variants.
```

## Remove chimeras

The core dada method removes substitution and indel errors, but chimeras remain. Fortunately, the accuracy of the sequences after denoising makes identifying chimeras simpler than it is when dealing with fuzzy OTUs: all sequences which can be exactly reconstructed as a bimeras (two-parent chimera) from more abundant sequences.

```
# Remove chimeric sequences:
seqtab.nochim <- removeBimeraDenovo(seqtab, method = "consensus",
  multithread = threads, verbose = TRUE)
```

```
## Identified 934 bimeras out of 1452 input sequences.
```

```
dim(seqtab.nochim)

## [1] 6 518
# ratio of chimeric sequence reads
1 - sum(seqtab.nochim)/sum(seqtab)

## [1] 0.1049661
# write sequence variants count table to file
write.table(t(seqtab.nochim), paste0(output.dir, "/all_samples_SV-counts.tsv"),
  sep = "\t", eol = "\n", quote = F, col.names = NA)
# write OTU table to file
saveRDS(seqtab.nochim, paste0(output.dir, "/seqtab_final.rds"))
```

**IMPORTANT:** Most of your **reads** should remain after chimera removal (it is not uncommon for a majority of **sequence variants** to be removed though). If most of your reads were removed as chimeric, upstream processing may need to be revisited. In almost all cases this is caused by primer sequences with ambiguous nucleotides that were not removed prior to beginning the DADA2 pipeline.

## Track reads through the pipeline

As a final check of the progress, look at the number of reads that made it through each step in the pipeline. This is a great place to do a last sanity check. Outside of filtering (depending on how stringent you want to be) there should no step in which a majority of reads are lost. If a majority of reads failed to merge, you may need to revisit the truncLen parameter used in the filtering step and make sure that the truncated reads span your amplicon. If a majority of reads failed to pass the chimera check, you may need to revisit the removal of primers, as the ambiguous nucleotides in unremoved primers interfere with chimera identification.

```
getN <- function(x) sum(getUniques(x))
track <- cbind(rd.counts, sapply(mergers, getN), rowSums(seqtab),
  rowSums(seqtab.nochim))
colnames(track) <- c("input", "filtered", "ratio", "merged",
  "tabled", "nonchim")
rownames(track) <- sample.names
# print table
track

##           input filtered ratio merged tabled nonchim
## Orwoll_BI0023_BI 62724    53296 0.85 49613 49613 45189
## Orwoll_BI0056_BI 55342    47545 0.86 45841 45841 39776
## Orwoll_BI0131_BI 55144    46797 0.85 44226 44226 42032
## Orwoll_BI0153_BI 44610    39510 0.89 37388 37388 33152
## Orwoll_BI0215_BI 48227    40256 0.83 37812 37812 33457
## Orwoll_BI0353_BI 54271    49323 0.91 46633 46633 40457

# save to file
write.table(track, paste0(readQC.folder, "/Read_counts_at_each_step.tsv"),
  sep = "\t", quote = F, eol = "\n", col.names = NA)

# rename fasta header and write to file
seqtab.nochim_trans <- as.data.frame(t(seqtab.nochim)) %>%
  rownames_to_column(var = "sequence") %>%
  rowid_to_column(var = "OTUNumber") %>%
  mutate(OTUNumber = sprintf("otu%04d", OTUNumber)) %>%
  mutate(sequence = str_replace_all(sequence, "(-|\\.|\\.)", ""))
```

```
df <- seqtab.nochim_trans
seq_out <- Biostrings::DNAStringSet(df$sequence)

names(seq_out) <- df$OTUNumber

Biostrings::writeXStringSet(seq_out, str_c(output.dir, "/ASV_rep_no_tax.fasta"),
  compress = FALSE, width = 20000)
```

## Align sequences and reconstruct phylogeny

Multiple sequence alignment of resolved sequence variants is used to generate a phylogenetic tree, which is required for calculating UniFrac beta-diversity distances between microbiome samples.

```
# Get sequences
seqs <- getSequences(seqtab.nochim)
# save non-chimera, representative sequences

names(seqs) <- seqs # This propagates to the tip labels of the tree
# Multiple sequence alignment
mult <- msa(seqs, method = "ClustalOmega", type = "dna", order = "input")

## using Gonnet
# Save msa to file; convert first to phangorn object
phang.align <- as.phyDat(mult, type = "DNA", names = getSequences(seqtab.nochim))
write.phyDat(phang.align, format = "fasta", file = paste0(output.dir,
  "/msa.fasta"))

# Call FastTree (via 'system') to reconstruct phylogeny
if (unname(Sys.info()["sysname"]) == "Linux") {
  system(paste("../dependencies/FastTreeMP_Linux -gtr -nt ",
    output.dir, "/msa.fasta > ", output.dir, "/FastTree.tre",
    sep = ""))
} else {
  system(paste("../dependencies/FastTreeMP_Darwin -gtr -nt ",
    output.dir, "/msa.fasta > ", output.dir, "/FastTree.tre",
    sep = ""))
}

detach("package:phangorn", unload = TRUE)
detach("package:msa", unload = TRUE)
```

## Assign taxonomy

The assignTaxonomy function takes a set of sequences and a training set of taxonomically classified sequences, and outputs the taxonomic assignments with at least minBoot bootstrap confidence. Formatted training datasets for taxonomic assignments can be downloaded from here <https://benjjneb.github.io/dada2/training.html>.

assignTaxonomy( ... ) implements the RDP naive Bayesian classifier method described in Wang et al. 2007. In short, the kmer profile of the sequences to be classified are compared against the kmer profiles of all sequences in a training set of sequences with assigned taxonomies. The reference sequence with the most similar profile is used to assign taxonomy to the query sequence, and then a bootstrapping approach is used to assess the confidence assignment at each taxonomic level. The return value of assignTaxonomy(...) is a

character matrix, with each row corresponding to an input sequence, and each column corresponding to a taxonomic level.

```
# Assign taxonomy: print(paste('Assign taxonomy (Greengene)
# starts at', Sys.time, sep = ' '))
taxa.gg13_8 <- assignTaxonomy(seqtab.nochim, "../reference_dbs_16S/gg_13_8_train_set_97.fa.gz",
    multithread = threads, tryRC = TRUE)

# Print first 6 rows of taxonomic assignment
unname(head(taxa.gg13_8))

##      [,1]      [,2]      [,3]
## [1,] "k__Bacteria" "p__Bacteroidetes" "c__Bacteroidia"
## [2,] "k__Bacteria" "p__Proteobacteria" "c__Gammaproteobacteria"
## [3,] "k__Bacteria" "p__Bacteroidetes" "c__Bacteroidia"
## [4,] "k__Bacteria" "p__Firmicutes" "c__Clostridia"
## [5,] "k__Bacteria" "p__Bacteroidetes" "c__Bacteroidia"
## [6,] "k__Bacteria" "p__Bacteroidetes" "c__Bacteroidia"
##      [,4]      [,5]      [,6]
## [1,] "o__Bacteroidales" "f__Bacteroidaceae" "g__Bacteroides"
## [2,] "o__Enterobacteriales" "f__Enterobacteriaceae" "g__Klebsiella"
## [3,] "o__Bacteroidales" "f__Bacteroidaceae" "g__Bacteroides"
## [4,] "o__Clostridiales" "f__Ruminococcaceae" "g__Faecalibacterium"
## [5,] "o__Bacteroidales" "f__Bacteroidaceae" "g__Bacteroides"
## [6,] "o__Bacteroidales" "f__Bacteroidaceae" "g__Bacteroides"
##      [,7]
## [1,] "s__"
## [2,] "s__"
## [3,] "s__caccae"
## [4,] "s__prausnitzii"
## [5,] "s__ovatus"
## [6,] "s__ovatus"

# Replace NAs in taxonomy assignment table with prefix
# corresponding to tax rank
taxa.gg13_8.2 <- replaceNA.in.assignedTaxonomy(taxa.gg13_8)

# Write taxa table to file
write.table(taxa.gg13_8.2, paste0(output.dir, "/all_samples_GG13-8-taxonomy.tsv"),
    sep = "\t", eol = "\n", quote = F, col.names = NA)
```

## Merge OTU and GG13-8 taxonomy tables

```
otu.gg.tax.table <- merge(t(seqtab.nochim), taxa.gg13_8.2, by = "row.names")
rownames(otu.gg.tax.table) <- otu.gg.tax.table[, 1]
otu.gg.tax.table <- otu.gg.tax.table[, -1]

write.table(otu.gg.tax.table, paste0(output.dir, "/all_samples_SV-counts_and_GG13-8-taxonomy.tsv"),
    sep = "\t", eol = "\n", quote = F, col.names = NA)
## print(paste('Assign taxonomy (Greengene) finishes
## at', Sys.time, sep = ' '))
```

For RDP and Silva, taxonomic assignment to species level is a two-step process. Fast and appropriate species-level assignment from 16S data is provided by the `assignSpecies( ... )` method. `assignSpecies( ... )` uses exact string matching against a reference database to assign Genus species binomials. In short,

query sequence are compared against all reference sequences that had binomial genus-species nomenclature assigned, and the genus-species of all exact matches are recorded and returned if it is unambiguous.

The convenience function `addSpecies( ... )` takes as input a taxonomy table, and outputs a table with an added species column. Only those genus-species binomials which are consistent with the genus assigned in the provided taxonomy table are retained in the output. See [here](https://benjjneb.github.io/dada2/assign.html) for more on taxonomic assignment <https://benjjneb.github.io/dada2/assign.html>.

## Assign SILVA and RDP taxonomies and merge with OTU table

```
# Assign SILVA taxonomy print(paste('Assign taxonomy (Silva)
# starts at', Sys.time, sep = ' '))
taxa.silva <- assignTaxonomy(seqtab.nochim, "../reference_dbs_16S/silva_nr_v128_train_set.fa.gz",
  multithread = threads)
saveRDS(taxa.silva, paste0(output.dir, "/taxa.silva.1.rds"))
# Replace NAs in taxonomy assignment table with prefix
# corresponding to tax rank
taxa.silva.2 <- replaceNA.in.assignedTaxonomy(taxa.silva)
saveRDS(taxa.silva.2, paste0(output.dir, "/taxa.silva.2.rds"))
# OMIT APPENDING SPECIES FOR SILVA DUE TO MEMORY CONSTRAINTS
# Append species. Note that appending the argument
# 'allowMultiple=3' will return up to 3 different matched
# species, but if 4 or more are matched it returns NA.
# taxa.silva.species <- addSpecies(taxa.silva,
# '/n/huttenhower_lab/data/dada2_reference_databases/silva_species_assignment_v128.fa.gz')

# Merge with OTU table and save to file
otu.silva.tax.table <- merge(t(seqtab.nochim), taxa.silva.2,
  by = "row.names")
rownames(otu.silva.tax.table) <- otu.silva.tax.table[, 1]
otu.silva.tax.table <- otu.silva.tax.table[, -1]

write.table(otu.silva.tax.table, paste0(output.dir, "/all_samples_SV-counts_and_SILVA-taxonomy.tsv"),
  sep = "\t", eol = "\n", quote = F, col.names = NA)
## print(paste('Assign taxonomy (Silva) finishes at', Sys.time,
## sep = ' ')) Assign RDP taxonomy print(paste('Assign
## taxonomy (RDP) starts at', Sys.time, sep = ' '))
taxa.rdp <- assignTaxonomy(seqtab.nochim, "../reference_dbs_16S/rdp_train_set_18.fa.gz",
  multithread = threads)

## Warning in .Call2("fasta_index", filexp_list, nrec, skip, seek.first.rec, :
## reading FASTA file ../reference_dbs_16S/rdp_train_set_18.fa.gz: ignored 9
## invalid one-letter sequence codes

# Replace NAs in taxonomy assignment table with prefix
# corresponding to tax rank
taxa.rdp.2 <- replaceNA.in.assignedTaxonomy(taxa.rdp)

# OMIT APPENDING SPECIES FOR RDP DUE TO MEMORY CONSTRAINTS
# Append species. Note that appending the argument
# 'allowMultiple=3' will return up to 3 different matched
# species, but if 4 or more are matched it returns NA.
# taxa.rdp.species <- addSpecies(taxa.rdp,
# '/n/huttenhower_lab/data/dada2_reference_databases/rdp_species_assignment_16.fa.gz')
```

```

# Merge with OTU table and save to file
otu.rdp.tax.table <- merge(t(seqtab.nochim), taxa.rdp.2, by = "row.names")
rownames(otu.rdp.tax.table) <- otu.rdp.tax.table[, 1]
otu.rdp.tax.table <- otu.rdp.tax.table[, -1]
# extract representative sequences and save to file

write.table(otu.rdp.tax.table, paste0(output.dir, "/all_samples_SV-counts_and_RDP-taxonomy.tsv"),
  sep = "\t", eol = "\n", quote = F, col.names = NA)
## print(paste('Assign taxonomy (RDP) finishes at', Sys.time,
## sep = ' '))

```

Session Info:

```
sessionInfo()
```

```
## R version 4.0.2 (2020-06-22)
## Platform: x86_64-conda_cos6-linux-gnu (64-bit)
## Running under: Red Hat Enterprise Linux
##
## Matrix products: default
## BLAS/LAPACK: /storage/coda1/p-ktk3/0/jzhao399/miniconda3/lib/libopenblas-r0.3.12.so
##
## locale:
##  [1] LC_CTYPE=en_US.UTF-8      LC_NUMERIC=C
##  [3] LC_TIME=en_US.UTF-8      LC_COLLATE=en_US.UTF-8
##  [5] LC_MONETARY=en_US.UTF-8  LC_MESSAGES=en_US.UTF-8
##  [7] LC_PAPER=en_US.UTF-8     LC_NAME=C
##  [9] LC_ADDRESS=C             LC_TELEPHONE=C
## [11] LC_MEASUREMENT=en_US.UTF-8 LC_IDENTIFICATION=C
##
## attached base packages:
## [1] stats4      parallel  stats      graphics  grDevices  utils      datasets
## [8] methods    base
##
## other attached packages:
##  [1] stringr_1.4.0      RcppParallel_5.0.3
##  [3] knitr_1.31         rmarkdown_2.7
##  [5] gridExtra_2.3      ape_5.4-1
##  [7] dada2_1.18.0       Rcpp_1.0.5
##  [9] Rsamtools_2.6.0    Biocstrings_2.58.0
## [11] XVector_0.30.0     BiocParallel_1.24.1
## [13] SummarizedExperiment_1.20.0 Biobase_2.50.0
## [15] MatrixGenerics_1.2.1 matrixStats_0.58.0
## [17] GenomicRanges_1.42.0 GenomeInfoDb_1.26.7
## [19] IRanges_2.24.1     S4Vectors_0.28.1
## [21] BiocGenerics_0.36.1 tibble_3.1.0
## [23] cowplot_1.1.1      ggplot2_3.3.3
## [25] mgcv_1.8-34        nlme_3.1-150
## [27] RCurl_1.98-1.3     devtools_2.3.2
## [29] usethis_2.0.1      dplyr_1.0.2
## [31] BiocManager_1.30.12
##
## loaded via a namespace (and not attached):
##  [1] bitops_1.0-7      fs_1.5.0          RColorBrewer_1.1-2
##  [4] rprojroot_2.0.2   tools_4.0.2       utf8_1.2.1
##  [7] R6_2.5.0          colorspace_2.0-0  withr_2.4.2
## [10] tidyselect_1.1.0  prettyunits_1.1.1 processx_3.5.1
## [13] compiler_4.0.2    cli_2.5.0         formatR_1.9
## [16] desc_1.3.0        DelayedArray_0.16.3 labeling_0.4.2
## [19] scales_1.1.1      quadprog_1.5-8    callr_3.6.0
## [22] digest_0.6.27     htmltools_0.5.0   jpeg_0.1-8.1
## [25] pkgconfig_2.0.3   sessioninfo_1.1.1 highr_0.8
## [28] fastmap_1.0.1     rlang_0.4.10      rstudioapi_0.13
## [31] farver_2.1.0      generics_0.1.0    hwriter_1.3.2
## [34] magrittr_2.0.1    GenomeInfoDbData_1.2.4 Matrix_1.3-2
## [37] munsell_0.5.0     fansi_0.4.2       lifecycle_1.0.0
```



## [40] yaml_2.2.1	stringi_1.5.3	zlibbioc_1.36.0
## [43] pkgbuild_1.2.0	plyr_1.8.6	grid_4.0.2
## [46] crayon_1.4.1	lattice_0.20-41	splines_4.0.2
## [49] ps_1.6.0	pillar_1.5.1	igraph_1.2.6
## [52] reshape2_1.4.4	codetools_0.2-18	pkgload_1.1.0
## [55] fastmatch_1.1-0	glue_1.4.2	evaluate_0.14
## [58] ShortRead_1.48.0	latticeExtra_0.6-29	remotes_2.2.0
## [61] png_0.1-7	vctrs_0.3.7	testthat_3.0.0
## [64] gtable_0.3.0	purrr_0.3.4	cachem_1.0.4
## [67] xfun_0.20	GenomicAlignments_1.26.0	memoise_2.0.0
## [70] ellipsis_0.3.2		