



Data Analysis Report

GATC Microbiome Profiling (Combined Analysis) v3.6

Project(s): NG-13133

October 4, 2017

Table of Contents

1	Introduction	1
2	Samples	1
3	Analysis Summary	3
3.1	Workflow	3
3.2	Merging read pairs by overlapping	4
3.3	Clustering by sequence similarity	4
3.4	Chimera check and removal	4
3.5	OTU assignment	4
3.6	Filtering and plotting	5
3.7	Microbial Diversity Indices	5
3.7.1	Overview	5
3.7.2	Shannon Index	6
3.7.3	Simpson Index	6
4	Results	7
4.1	Read statistics	7
4.2	Diversity index tables	7
4.3	OTU abundance tables	7
4.4	Phylum	8
4.5	Class	11
4.6	Order	14
4.7	Family	17
4.8	Genus	20
4.9	Species	23
5	Deliverables	26
6	Formats	26
7	Tables	27
8	FAQ	28
	Bibliography	29

1 Introduction

The analysis of genes common to or ubiquitous amongst various organisms like bacterial 16S rRNA or fungal ITS is a time- and cost-effective method to characterise microbial diversity in complex samples.

Amplification and high-throughput sequencing of the hypervariable regions of these genes is therefore a commonly used method for studying phylogeny and taxonomy. It is particularly suitable for analysing diverse samples and unculturable microorganisms and is therefore usable for various industrial, agricultural, medical and environmental applications.

This amplicon-based method has been optimised regarding study design and bioinformatics processing to provide a ready to use solution for researchers who are seeking to characterise microbiomes from various sources and samples which are usually difficult to study.



Figure 1: Schematic overview of the 16S rRNA gene. The sequence identity of the 16S rRNA gene of more than 6,000 bacteria compared to consensus sequence is shown. Dips indicate hypervariable regions. Hypervariable regions (V1-V9) are shown in grey and the conserved regions in orange.

2 Samples

Table 1: 16S Primers used.

Variable Region	Primer	Sequence	Product size ¹
V3-V5[1]	357F	CCTACGGGAGGCAGCAG	570 bp
	926R	CCGTCAATTCMTTTRAGT	

Table 2: Analysed samples.

Sample	File Name
Project_1_B1_t144	NG-13133_Project_1_B1_t144_lib196227_5593_2_1.fastq
	NG-13133_Project_1_B1_t144_lib196227_5593_2_2.fastq
Project_1_B1_t48	NG-13133_Project_1_B1_t48_lib196224_5593_2_1.fastq
	NG-13133_Project_1_B1_t48_lib196224_5593_2_2.fastq

¹excluding primer lengths

Table 2: Analysed samples.

Sample	File Name
Project_1_G1_t144	NG-13133_Project_1_G1_t144.lib196228.5593_2_1.fastq
	NG-13133_Project_1_G1_t144.lib196228.5593_2_2.fastq
Project_1_G1_t48	NG-13133_Project_1_G1_t48.lib196225.5593_2_1.fastq
	NG-13133_Project_1_G1_t48.lib196225.5593_2_2.fastq
Project_1_G2_t144	NG-13133_Project_1_G2_t144.lib196229.5593_2_1.fastq
	NG-13133_Project_1_G2_t144.lib196229.5593_2_2.fastq
Project_1_G2_t48	NG-13133_Project_1_G2_t48.lib196226.5593_2_1.fastq
	NG-13133_Project_1_G2_t48.lib196226.5593_2_2.fastq
Project_1_t0	NG-13133_Project_1_t0.lib196223.5593_2_1.fastq
	NG-13133_Project_1_t0.lib196223.5593_2_2.fastq

3 Analysis Summary

3.1 Workflow

The schematic diagram of data analysis performed is displayed in the following graphic.

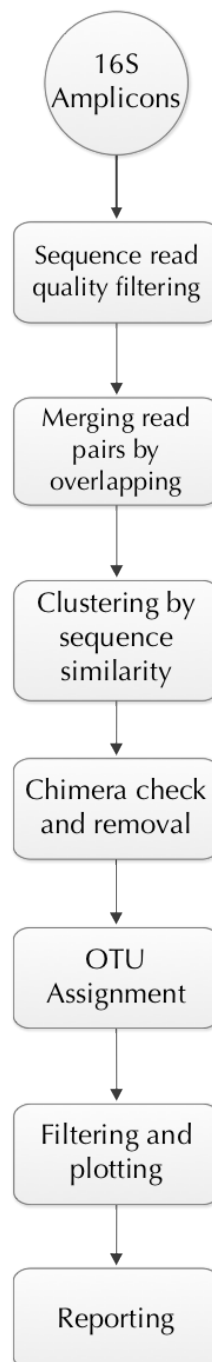


Figure 2: Microbiome Profiling Workflow

3.2 Merging read pairs by overlapping

In case of paired end sequencing where amplicons are sequenced in both the directions, the resulting read pairs are merged based on overlapping bases using FLASH[2] with maximum mismatch density of 0.25. Merging read pairs extends the read length to reflect the amplicon length which increases the possibility and accuracy of OTU assignment during the downstream processing.

3.3 Clustering by sequence similarity

In any given sample, there is an uneven representation of the microbiome biota which results in uneven amplification and sequence coverage. In order to reduce the computational time incurring for further downstream processing, the sequence data is compressed by performing sequence clustering based on 99% similarity accounting for PCR and sequencing errors (<1%). To achieve this, cd-hit[3], a clustering program is used. At high sequencing depth each original template is sequenced multiple times. Therefore singletons, clusters containing only one sequence, are removed from further analysis.

3.4 Chimera check and removal

PCR is an essential step in generating the amplicons from DNA samples. Due to the high similarity of different 16S rRNA, the possibility that small amounts of chimeric PCR products are generated is high. Therefore, the clustered data is checked for chimeras and the corresponding clusters are removed from further analysis. Chimera check is performed with UCHIME[4] using a full length, good quality, and non-chimeric 16S rRNA gene reference database.

3.5 OTU assignment

Non-chimeric, unique clusters are then subjected to BLASTn[5] analysis using non-redundant 16S rRNA reference sequences with an E-value cutoff of 1e-06. Reference 16S rRNA sequences are obtained from Ribosomal Database Project[6] (RDP Release 11 updated on September, 2016). Only good quality and unique 16S rRNA sequences which have a taxonomic assignment are considered and used as a reference database to assign operational taxonomic unit (OTU) status to the clusters. Taxonomic classification is based on NCBI Taxonomy[7] - <http://www.ncbi.nlm.nih.gov/taxonomy>. The number of sequences and length characteristics of the reference database used are described in table 3.

Table 3: Number of sequences and length characteristics for the reference database.

Total Sequences	Biggest	Smallest	Mean
11,795	1,768 bp	1,200 bp	1,461 bp

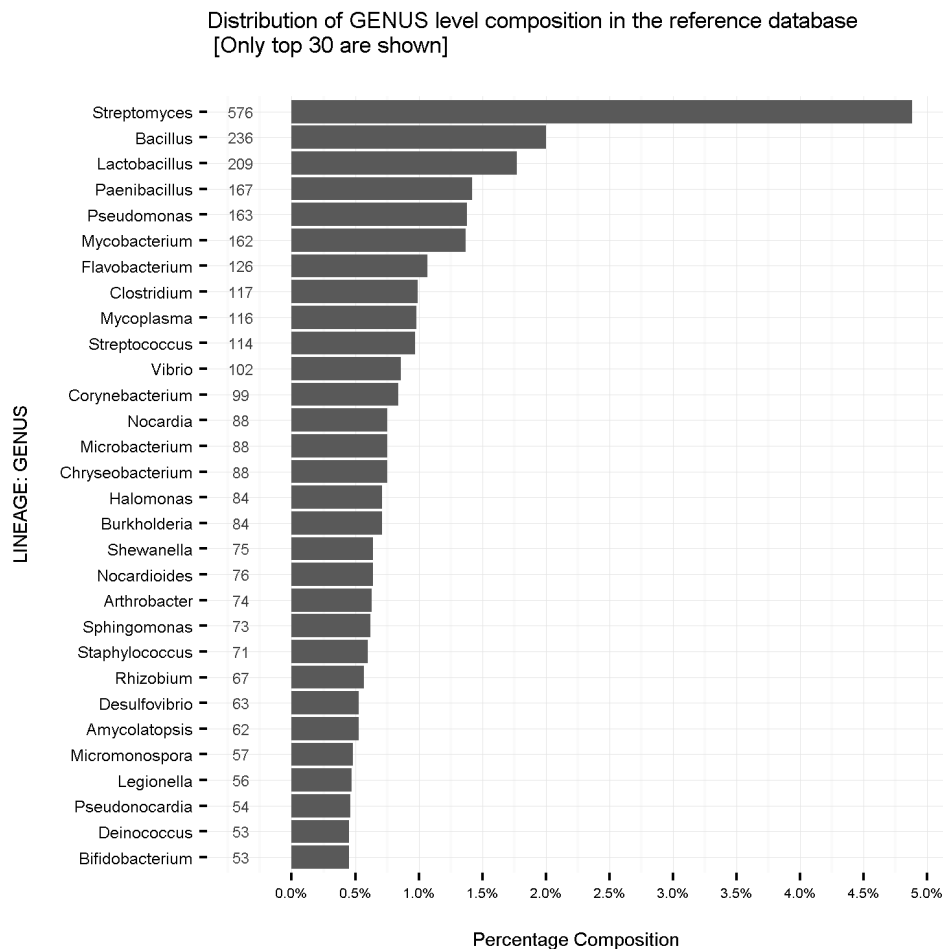


Figure 3: GENUS Distribution plot for the reference database.

3.6 Filtering and plotting

Except E-value cutoff ($1e-06$), no other thresholds were used during the BLAST analysis. All the hits to reference 16S rRNA database are considered and specific filters are applied to the hits to remove false positives. Further, best hit per cluster and multiple hits per cluster were analysed separately to determine the discriminatory power of the clusters with respect to the assigned OTUs. The various thresholds applied are reported in table 19. Finally, classification of OTU clusters and size of OTUs (number of reads within one cluster) are consolidated to compute relative abundancies (percentage composition).

3.7 Microbial Diversity Indices

3.7.1 Overview

A diversity index is a quantitative measure that reflects how many different types (such as species) there are in a dataset, and simultaneously takes into account how evenly the basic entities (such as individuals) are distributed among those types. The value of a diversity index increases both when the number of species increases and when evenness increases. For a given number of species, the value of a diversity index is maximized when all species are equally abundant.

Although there are many indices available to measure the species diversity, the most popular ones used in sequence based OTU profiling are - Shannon index and Simpson index.

3.7.2 Shannon Index

The idea behind this index is that diversity of a community is similar to the amount of information contained in a sampled environment. The Shannon index increases as both the richness and the evenness of the community increase. Since the index incorporates both components of biodiversity (richness and evenness) it provides a simplistic summary of species diversity. On the other hand, it makes it difficult to compare communities that differ greatly in richness. In order to overcome this limitation, a second index - Simpson index is used for comparative studies, combining a direct estimate of species richness (the total number of species in the community - Shannon Index) with some measure of dominance or evenness (Simpson index) [8].

3.7.3 Simpson Index

It measures the probability of any two individuals randomly selected from a sample will belong to the same species. Simpson index gives the probability of any two individuals drawn from noticeably large community belonging to same species. Simpson index can be used to measure the species evenness, richness and diversity [9].

For general information about various diversity indices see [10].

4 Results

4.1 Read statistics

Table 4: Read statistics and OTU assignment

Sample	Total (Read pairs)	Cleaned	Merged by overlapping	Clustered by similarity	Chimeric	High quality	OTU assigned	Filter passed OTUs (best hit only)
Project_1_B1.t144	3,298,680	1,053,233	984,872	652,381	224,732	427,649	427,649	398,853
Project_1_B1.t48	2,074,420	663,370	631,004	380,155	205,775	174,380	174,380	128,677
Project_1_G1.t144	2,739,959	676,892	625,255	375,955	113,215	262,740	262,740	246,591
Project_1_G1.t48	2,689,623	693,362	651,147	357,922	144,566	213,356	213,356	94,686
Project_1_G2.t144	2,172,180	645,765	607,912	375,941	160,755	215,186	215,186	196,521
Project_1_G2.t48	2,557,593	830,144	791,669	488,314	117,144	371,170	371,170	311,256
Project_1.t0	2,210,881	500,333	459,111	226,821	468	226,353	226,351	30,167

4.2 Diversity index tables

Both Shannon and Simpson indices are computed using the R package Vegan^[11]. The results are found in the diversity index tables (Taxa-level.diversity_index_tables.tsv) and a graphical representation of the indices are in diversity plots (Taxa-level.diversity.png).

4.3 OTU abundance tables

Abundance measured by the percentage of OTU assigned reads from various taxonomic level was computed. The measured abundance levels are in OTU distribution tables (Taxa-level.combined.table.tsv) and the bar plots representing the abundance levels at various taxonomic level are in OTU distribution plots (Taxa-level.OTU.distribution.combined.png).

4.4 Phylum

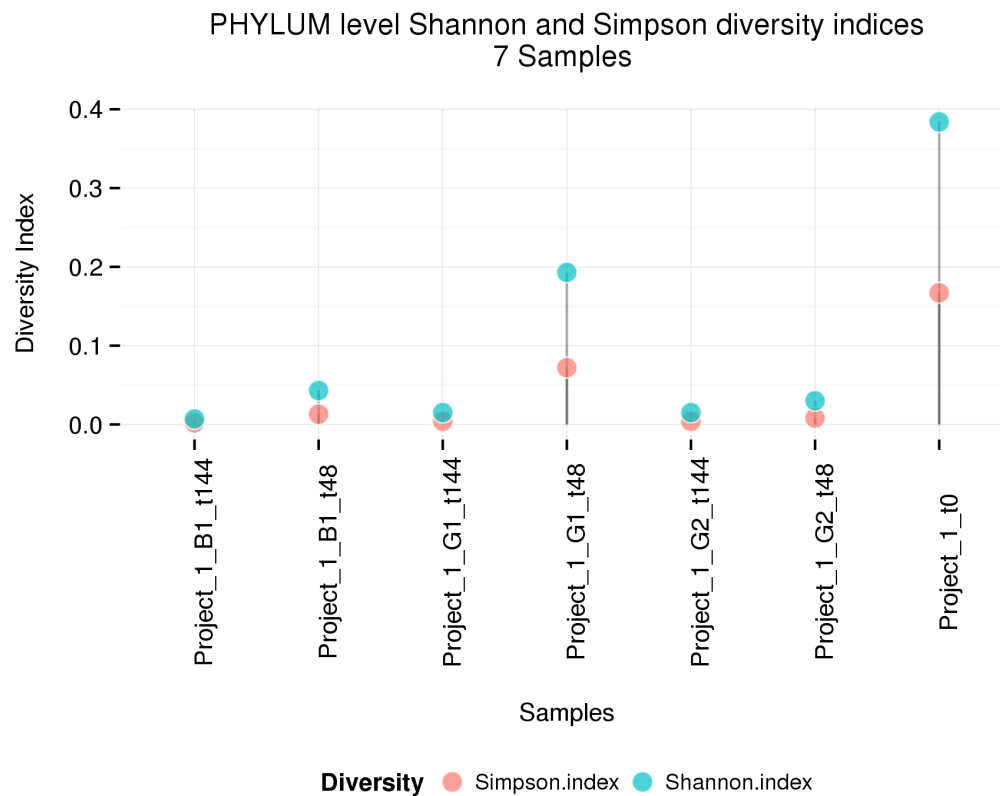


Figure 4: Phylum diversity indices (file: PHYLUM.diversity.png)

Table 5: Phylum diversity indices table (file: PHYLUM.diversity_index_tables.tsv)

Sample	Simpson.index	Shannon.index	OTUs
Project_1_B1_t144	0.002	0.007	2
Project_1_B1_t48	0.013	0.043	3
Project_1_G1_t144	0.004	0.015	3
Project_1_G1_t48	0.072	0.193	5
Project_1_G2_t144	0.004	0.015	3
Project_1_G2_t48	0.008	0.03	4
Project_1_t0	0.167	0.384	5

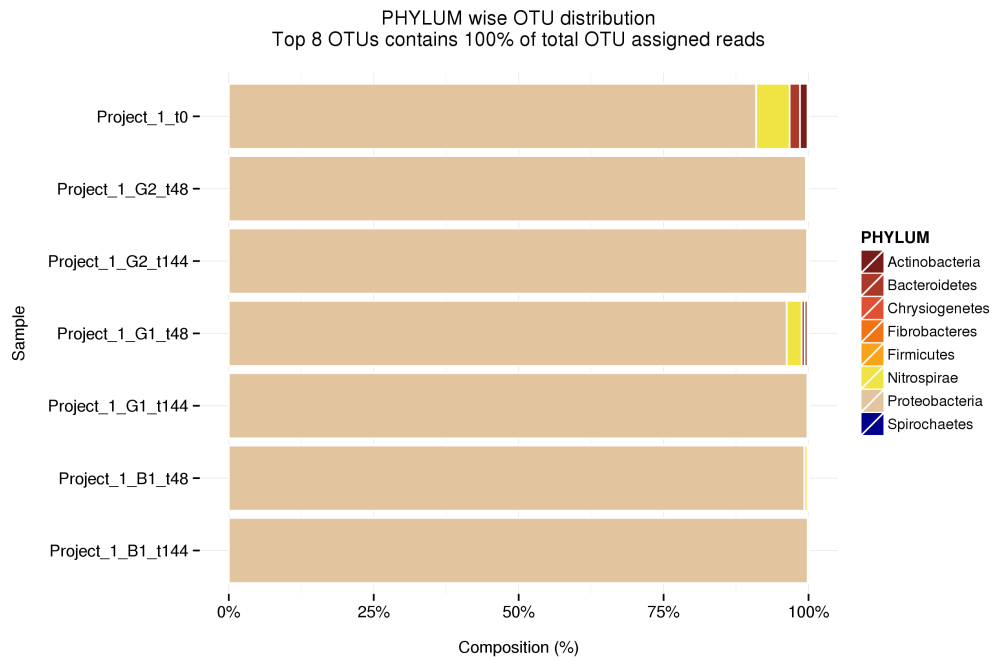


Figure 5: Phylum distribution plot (file: PHYLUM.OTU.distribution.combined.png)

Table 6: Phylum distribution table (file: PHYLUM.OTU.combined.table.percent.top.8.tsv)

PHYLUM	Project_1-B1-t144	Project_1-B1-t48	Project_1-G1-t144	Project_1-G1-t48	Project_1-G2-t144	Project_1-G2-t48	Project_1-t0
Proteobacteria	99.86	99.3	99.81	96.24	99.75	99.54	90.98
Nitrospirae	0.08	0.52	0.07	2.58	0.1	0.27	5.77
Bacteroidetes	0.05	0.14	0.12	0.53	0.09	0.07	1.77
Actinobacteria	0.01	0.04	0.01	0.48	0.01	0.02	1.33
Firmicutes	0	0	0	0.13	0.05	0.06	0.04
Spirochaetes	0	0	0	0.02	0	0	0.07
Fibrobacteres	0	0	0	0.03	0	0.01	0.04
Chrysiogenetes	0	0	0	0	0	0.04	0

4.5 Class

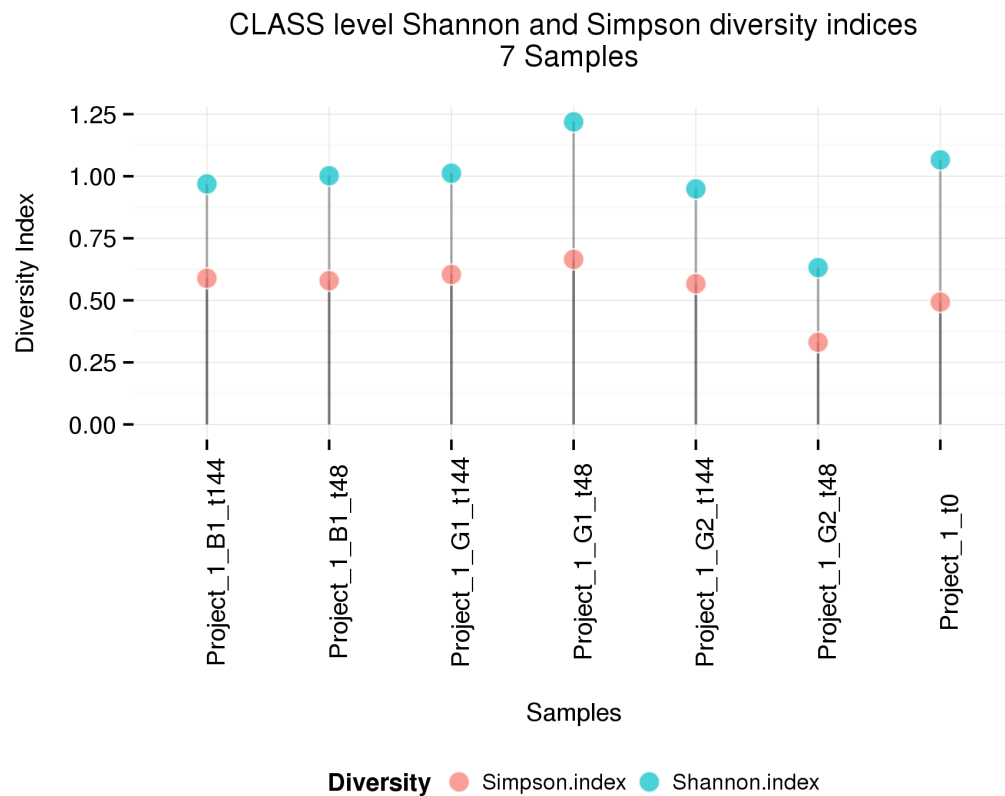


Figure 6: Class diversity indices (file: CLASS.diversity.png)

Table 7: Class diversity indices table (file: CLASS.diversity_index_tables.tsv)

Sample	Simpson.index	Shannon.index	OTUs
Project_1_B1_t144	0.589	0.969	4
Project_1_B1_t48	0.579	1.002	5
Project_1_G1_t144	0.604	1.012	6
Project_1_G1_t48	0.665	1.219	9
Project_1_G2_t144	0.567	0.949	4
Project_1_G2_t48	0.331	0.632	5
Project_1_t0	0.493	1.066	9

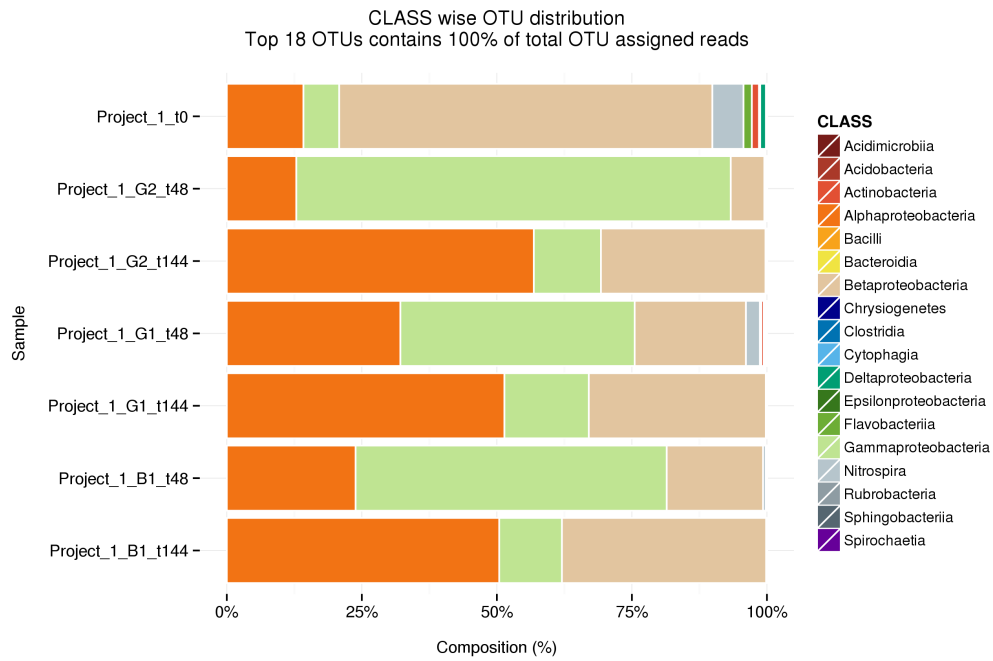


Figure 7: Class distribution plot (file: CLASS.OTU.distribution.combined.png)

Table 8: Class distribution table (file: CLASS.OTU.combined.table.percent.top.18.tsv)

CLASS	Project_1_B1_t144	Project_1_B1_t48	Project_1_G1_t144	Project_1_G1_t48	Project_1_G2_t144	Project_1_G2_t48	Project_1_t0
Alphaproteobacteria	50.44	23.84	51.4	32.14	56.85	12.87	14.19
Gammaaproteobacteria	11.61	57.61	15.62	43.37	12.41	80.43	6.62
Betaproteobacteria	37.8	17.83	32.77	20.6	30.47	6.23	69.08
Nitrospira	0.08	0.52	0.07	2.58	0.1	0.27	5.77
Flavobacteriia	0.03	0.12	0.06	0.27	0.05	0.06	1.55
Actinobacteria	0.01	0.04	0.01	0.46	0.01	0.02	1.32
Deltaproteobacteria	0.01	0.02	0.01	0.13	0.01	0.01	1.08
Sphingobacteriia	0.02	0.02	0.06	0.25	0.03	0.01	0.22
Clostridia	0	0	0	0.11	0.03	0.04	0.01
Spirochaetia	0	0	0	0.02	0	0	0.07
Acidobacteria	0	0	0	0.03	0	0.01	0.04
Bacilli	0	0	0	0.02	0.02	0.01	0.02
Chrysiogenetes	0	0	0	0	0	0.04	0
Rubrobacteria	0	0	0	0.01	0	0	0.01
Bacteroidia	0	0	0	0.01	0	0	0
Acidimicrobiia	0	0	0	0.01	0	0	0
Cytophagia	0	0	0	0	0	0	0
Epsilonproteobacteria	0	0	0	0	0	0	0

4.6 Order

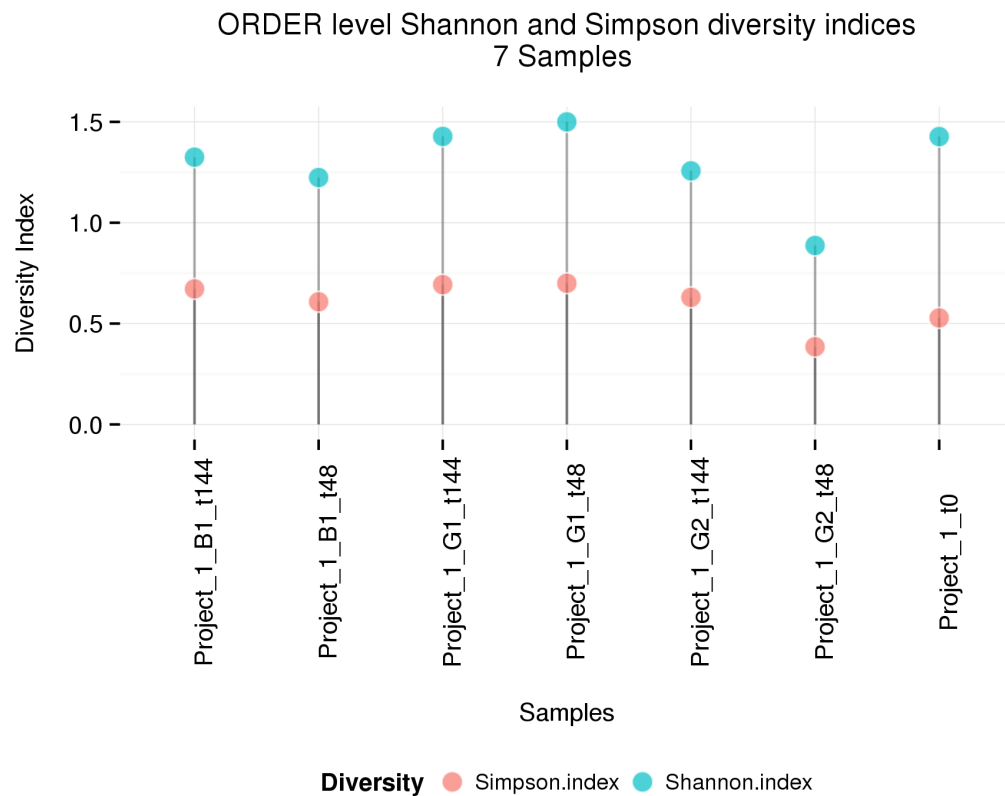


Figure 8: Order diversity indices (file: ORDER.diversity.png)

Table 9: Order diversity indices table (file: ORDER.diversity_index_tables.tsv)

Sample	Simpson.index	Shannon.index	OTUs
Project_1_B1_t144	0.672	1.325	9
Project_1_B1_t48	0.608	1.224	11
Project_1_G1_t144	0.693	1.428	14
Project_1_G1_t48	0.7	1.5	23
Project_1_G2_t144	0.63	1.257	10
Project_1_G2_t48	0.385	0.887	15
Project_1_t0	0.528	1.427	26

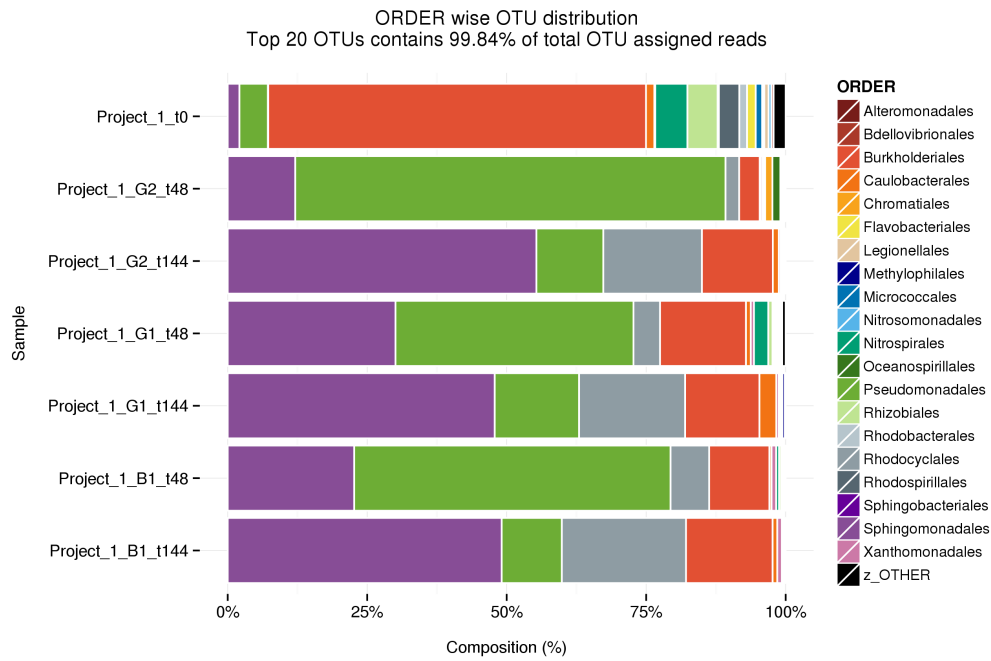


Figure 9: Order distribution plot (file: ORDER.OTU.distribution.combined.png)

Table 10: Order distribution table (file: ORDER.OTU.combined.table.percent.top.20.tsv)

ORDER	Project_1_B1_t144	Project_1_B1_t48	Project_1_G1_t144	Project_1_G1_t48	Project_1_G2_t144	Project_1_G2_t48	Project_1_t0
Pseudomonadales	10.78	56.7	15.12	42.65	11.99	77.11	5.09
Sphingomonadales	49.1	22.68	47.85	30.08	55.33	12.11	2.11
Burkholderiales	15.52	10.79	13.32	15.4	12.7	3.68	67.7
Rhodocyclales	22.25	6.92	19.01	4.77	17.68	2.44	0.07
Nitrospirales	0.08	0.52	0.07	2.58	0.1	0.27	5.77
Caulobacteriales	0.85	0.38	3.02	0.86	1.08	0.19	1.5
Rhizobiales	0.29	0.39	0.27	0.74	0.23	0.17	5.4
Rhodospirillales	0.15	0.3	0.17	0.3	0.16	0.15	3.67
Xanthomonadales	0.81	0.83	0.44	0.57	0.35	0.34	0.17
Flavobacteriales	0.03	0.12	0.06	0.27	0.05	0.06	1.55
Rhodobacteriales	0.04	0.07	0.09	0.17	0.05	0.25	1.37
Chromatiales	0.01	0.03	0.06	0.02	0.05	1.3	0.22
Micrococcales	0.01	0.03	0.01	0.3	0.01	0.01	1.16
Oceanospirillales	0	0	0	0	0.01	1.41	0.01
Legionellales	0.01	0.03	0.01	0.06	0.01	0.02	0.76
Methylophilales	0.03	0.03	0.44	0.13	0.08	0.05	0.02
Myxococcales	0	0	0	0.01	0	0	0.7
Nitrosomonadales	0	0.04	0	0.12	0	0.01	0.52
Sphingobacteriales	0.02	0.02	0.06	0.25	0.03	0.01	0.22
Gallionellales	0	0.02	0	0.08	0	0.02	0.45

4.7 Family

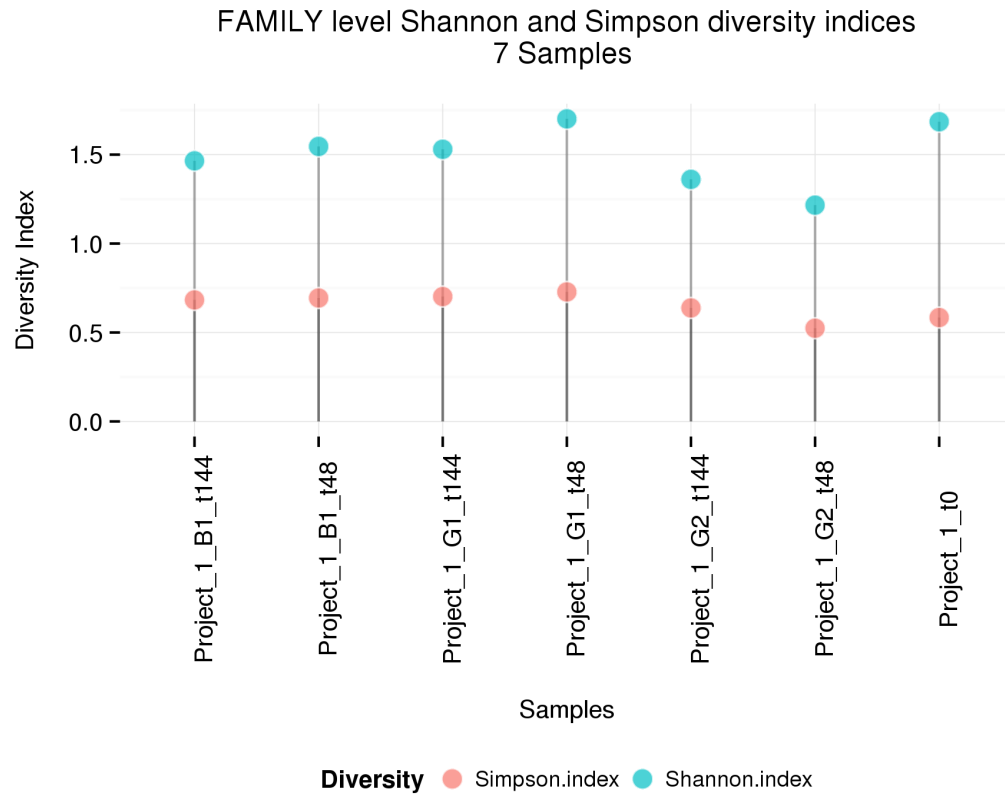


Figure 10: Family diversity indices (file: FAMILY.diversity.png)

Table 11: Family diversity indices table (file: FAMILY.diversity_index_tables.tsv)

Sample	Simpson.index	Shannon.index	OTUs
Project_1_B1_t144	0.683	1.465	13
Project_1_B1_t48	0.694	1.546	16
Project_1_G1_t144	0.702	1.53	16
Project_1_G1_t48	0.728	1.701	28
Project_1_G2_t144	0.638	1.361	13
Project_1_G2_t48	0.524	1.216	19
Project_1_t0	0.584	1.685	38

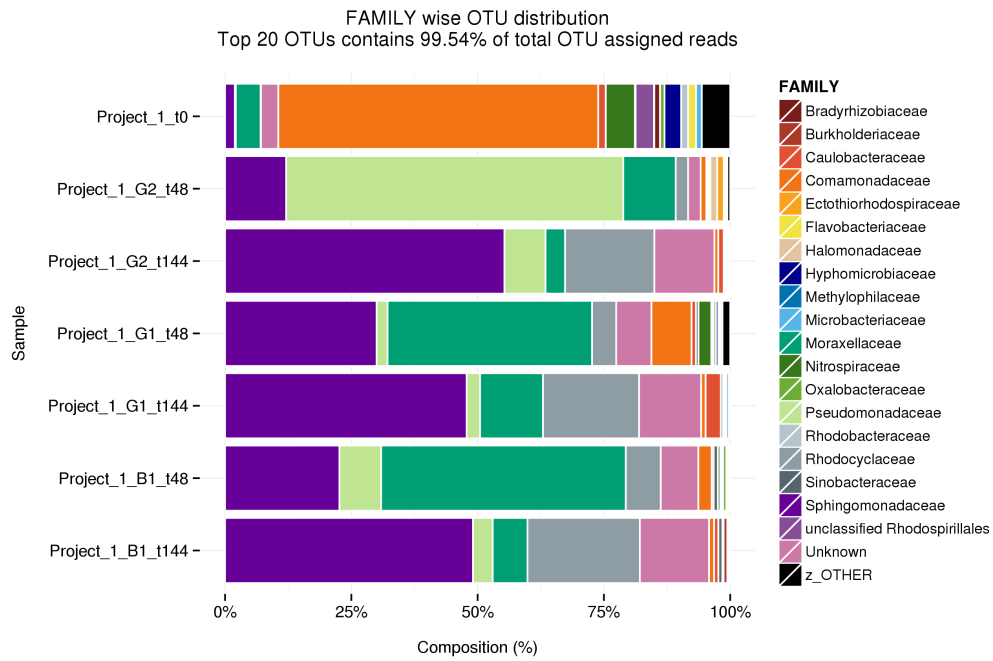


Figure 11: Family distribution plot (file: FAMILY.OTU.distribution.combined.png)

Table 12: Family distribution table (file: FAMILY.OTU.combined.table.percent.top.20.tsv)

FAMILY	Project_1_B1_t144	Project_1_B1_t48	Project_1_G1_t144	Project_1_G1_t48	Project_1_G2_t144	Project_1_G2_t48	Project_1_t0
Sphingomonadaceae	49.08	22.64	47.83	30.01	55.32	12.1	1.95
Moraxellaceae	6.88	48.43	12.46	40.49	3.87	10.38	4.91
Pseudomonadaceae	3.9	8.26	2.66	2.16	8.12	66.73	0.18
Comamonadaceae	1	2.65	0.88	7.93	0.75	1.13	63.32
Rhodocyclaceae	22.25	6.92	19.01	4.77	17.68	2.44	0.07
Unknown	13.7	7.44	12.3	6.98	11.9	2.49	3.44
Nitrospiraceae	0.08	0.52	0.07	2.58	0.1	0.27	5.77
Caulobacteraceae	0.85	0.38	3.02	0.86	1.08	0.19	1.5
unclassified Rhodospirillales	0.15	0.3	0.17	0.29	0.16	0.15	3.64
Hyphomicrobiaceae	0.04	0.19	0.05	0.47	0.06	0.08	3.31
Sinobacteraceae	0.8	0.82	0.43	0.52	0.34	0.33	0.01
Oxalobacteraceae	0.08	0.64	0.12	0.46	0.05	0.04	0.88
Flavobacteriaceae	0.03	0.12	0.06	0.27	0.05	0.06	1.55
Bradyrhizobiaceae	0.24	0.14	0.22	0.12	0.17	0.06	1.14
Rhodobacteraceae	0.04	0.07	0.09	0.17	0.05	0.25	1.37
Microbacteriaceae	0.01	0.03	0.01	0.26	0.01	0.01	1.15
Ectothiorhodospiraceae	0	0.03	0	0	0.01	1.3	0.12
Halomonadaceae	0	0	0	0	0	1.35	0.01
Legionellaceae	0.01	0.03	0.01	0.06	0.01	0.02	0.76
Burkholderiaceae	0.74	0.06	0.02	0.01	0	0	0.06

4.8 Genus

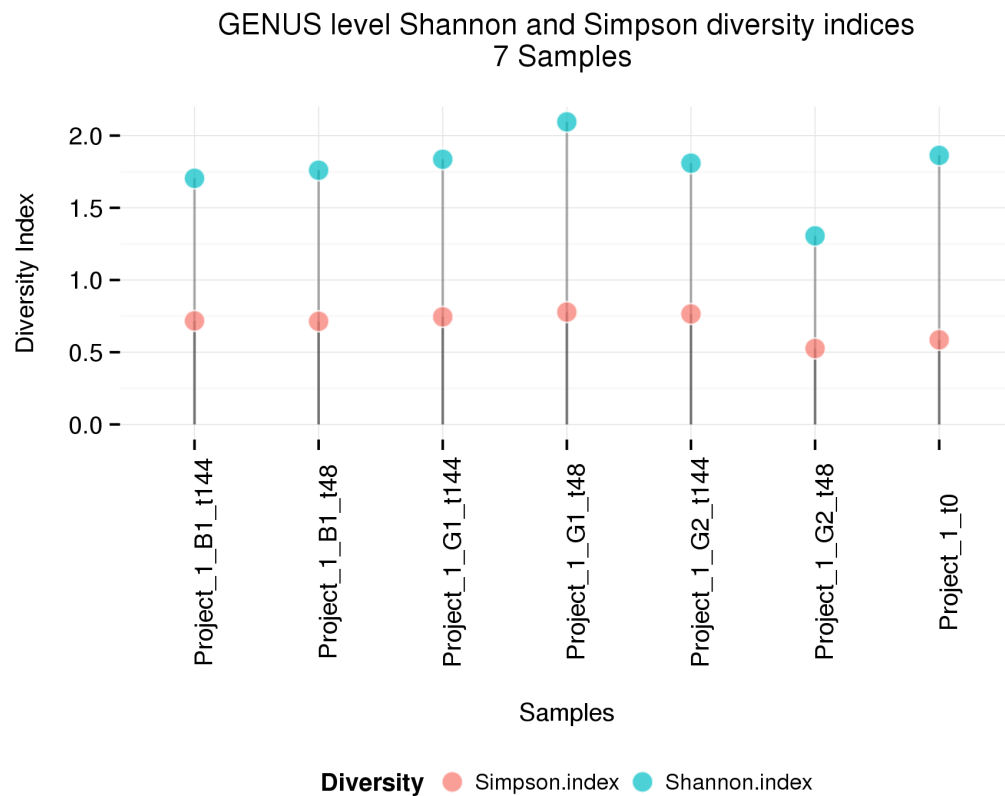


Figure 12: Genus diversity indices (file: GENUS.diversity.png)

Table 13: Genus diversity indices table (file: GENUS.diversity_index_tables.tsv)

Sample	Simpson.index	Shannon.index	OTUs
Project_1_B1_t144	0.717	1.704	24
Project_1_B1_t48	0.714	1.76	26
Project_1_G1_t144	0.745	1.836	30
Project_1_G1_t48	0.778	2.095	47
Project_1_G2_t144	0.765	1.809	26
Project_1_G2_t48	0.526	1.306	24
Project_1_t0	0.586	1.863	72

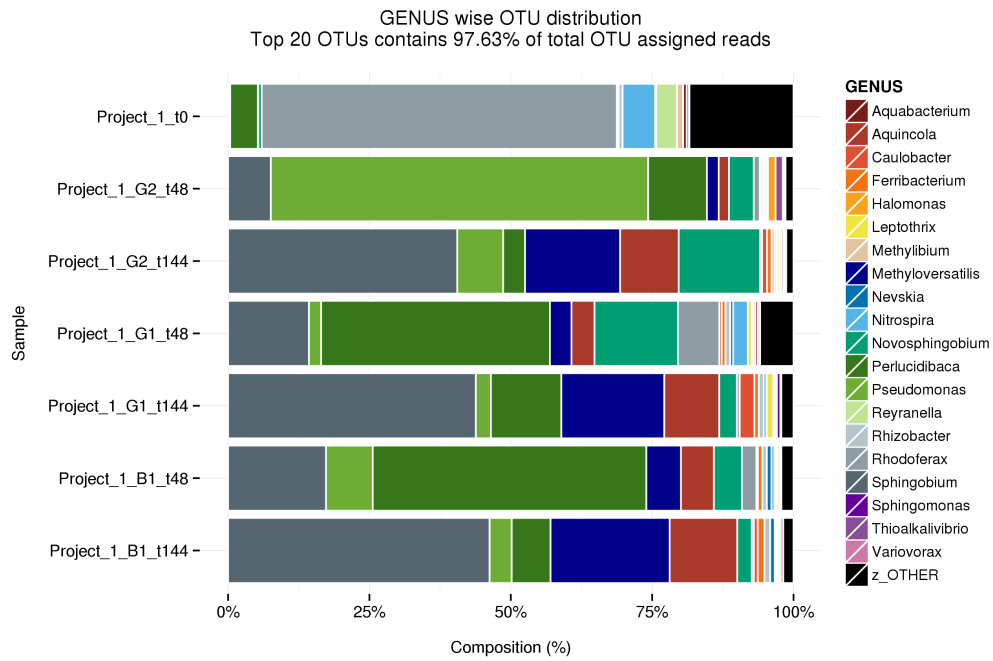


Figure 13: Genus distribution plot (file: GENUS.OTU.distribution.combined.png)

Table 14: Genus distribution table (file: GENUS.OTU.combined.table.percent.top.20.tsv)

GENUS	Project_1_B1_t144	Project_1_B1_t48	Project_1_G1_t144	Project_1_G1_t48	Project_1_G2_t144	Project_1_G2_t48	Project_1_t0
Sphingobium	46.25	17.32	43.82	14.31	40.53	7.57	0.22
Perlucidibaca	6.88	48.39	12.45	40.45	3.87	10.37	4.89
Pseudomonas	3.9	8.26	2.66	2.16	8.12	66.73	0.17
Rhodoferrax	0.38	2.55	0.51	7.32	0.33	1.06	62.8
Methyloversatilis	21.07	6.12	18.21	3.81	16.81	2.09	0.05
Novosphingobium	2.56	4.97	3.12	14.76	14.42	4.39	0.62
Aquicola	11.96	5.82	9.71	4.08	10.36	1.8	0.02
Nitrospira	0.08	0.52	0.07	2.58	0.1	0.27	5.77
Caulobacter	0.68	0.27	2.6	0.42	0.88	0.12	0.31
Rhizobacter	1.02	0.78	0.88	0.85	0.52	0.32	0.69
Reyranella	0.15	0.3	0.17	0.29	0.16	0.15	3.64
Ferribacterium	1.19	0.8	0.78	0.67	0.8	0.32	0
Hyphomicrobium	0.04	0.16	0.04	0.4	0.05	0.06	2.71
Nevskia	0.8	0.82	0.43	0.52	0.34	0.33	0.01
Leptothrix	0.23	0.17	1.19	0.69	0.4	0.09	0.08
Methylibium	0.25	0.39	0.16	0.39	0.19	0.14	1.09
Aquabacterium	0.2	0.21	0.29	0.44	0.34	0.09	0.63
Flavobacterium	0.03	0.12	0.06	0.26	0.05	0.06	1.53
Sphingomonas	0.14	0.11	0.6	0.34	0.23	0.05	0.41
Bradyrhizobium	0.18	0.11	0.13	0.1	0.12	0.05	1

4.9 Species

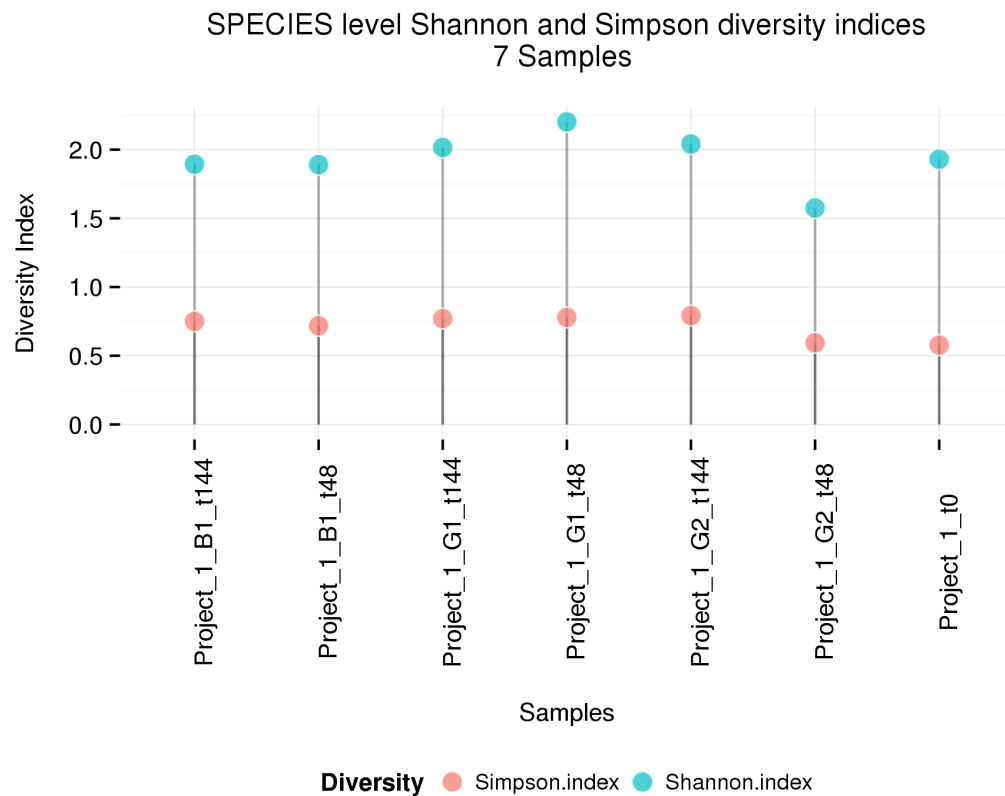


Figure 14: Species diversity indices (file: SPECIES.diversity.png)

Table 15: Species diversity indices table (file: SPECIES.diversity_index_tables.tsv)

Sample	Simpson.index	Shannon.index	OTUs
Project_1_B1_t144	0.75	1.894	35
Project_1_B1_t48	0.718	1.891	38
Project_1_G1_t144	0.77	2.015	42
Project_1_G1_t48	0.779	2.202	68
Project_1_G2_t144	0.792	2.041	39
Project_1_G2_t48	0.594	1.575	31
Project_1_t0	0.577	1.93	99

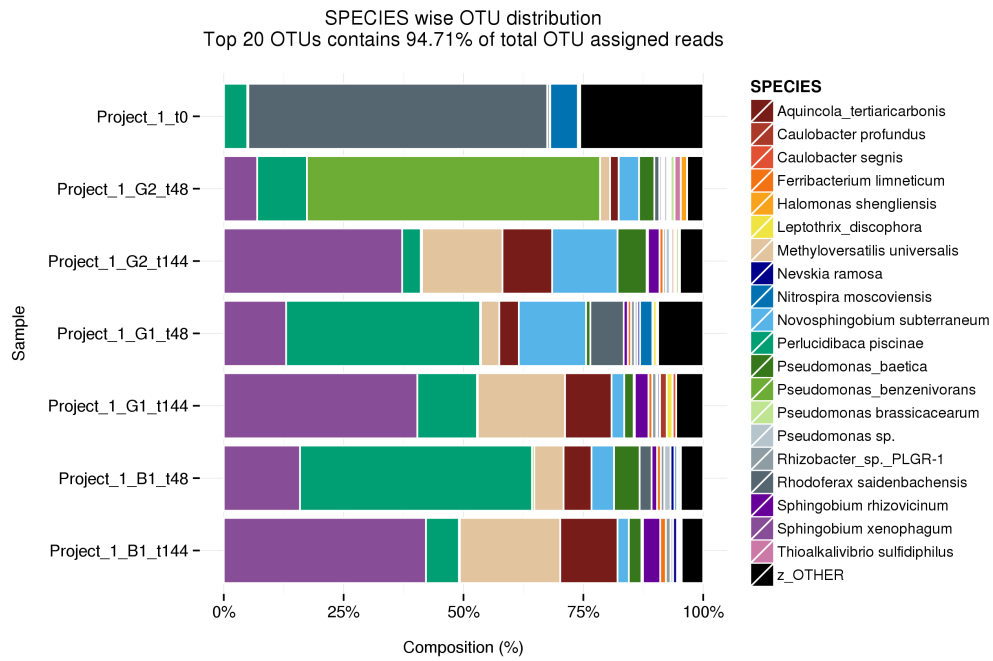


Figure 15: Species distribution plot (file: SPECIES.OTU.distribution.combined.png)

Table 16: Species distribution table (file: SPECIES.OTU.combined.table.percent.top.20.tsv)

SPECIES	Project_1_B1_t144	Project_1_B1_t48	Project_1_G1_t144	Project_1_G1_t48	Project_1_G2_t144	Project_1_G2_t48	Project_1_t0
Sphingobium xenophagum	42.18	15.91	40.38	13.01	37.21	6.98	0.04
Perleidibaca piscinae	6.88	48.39	12.45	40.45	3.87	10.37	4.89
Rhodofarax saidenbachensis	0.32	2.5	0.24	6.95	0.28	1.03	62.38
Methyloversatilis universalis	20.94	6.11	18.19	3.8	16.78	2.09	0.05
Pseudomonas_benzenivorans	0.19	0.47	0.16	0.2	0.27	61.13	0
Aquicola_tertiaricarbonis	11.96	5.82	9.71	4.08	10.36	1.8	0.02
Novosphingobium subterraneum	2.35	4.7	2.64	14.02	13.64	4.24	0.12
Pseudomonas_baetica	2.64	5.34	1.97	0.9	6.07	3.19	0
Sphingobium rhizoviciunum	3.59	1.14	2.85	0.87	2.42	0.41	0
Nitrospira moscoviensis	0.08	0.52	0.07	2.58	0.1	0.27	5.77
Rhizobacter_sp._PLGR-1	0.98	0.72	0.85	0.77	0.51	0.31	0.54
Ferribacterium limneticum	1.19	0.8	0.78	0.67	0.8	0.32	0
Pseudomonas sp.	0.53	1.3	0.27	0.6	0.82	0.6	0.03
Reyranella soli	0.13	0.25	0.14	0.23	0.13	0.12	2.92
Nevskia ramosa	0.8	0.82	0.43	0.52	0.34	0.33	0.01
Leptothrix discophora	0.23	0.17	1.19	0.69	0.4	0.09	0.08
Caulobacter profundus	0.28	0.13	1.4	0.18	0.45	0.06	0.21
Methylibium_petroleiphilum_PM1	0.25	0.39	0.16	0.39	0.19	0.14	1.09
Pseudomonas brassicacearum	0.13	0.43	0.06	0.21	0.59	0.8	0
Hyphomicrobium vulgare	0.02	0.08	0.02	0.28	0.01	0.04	1.57

5 Deliverables

Table 17: List of deliverable files, format and recommended programs to access.

File	Format	Program To Open File
Taxa-level.diversity_index_tables.tsv	TSV	Spreadsheet editor
Taxa-level.diversity.png	PNG	Image viewer
Taxa-level.combined.table.percent.top.X.tsv	TSV	Spreadsheet editor
Taxa-level.combined.table.percent.tsv	TSV	Spreadsheet editor
Taxa-level.combined.table.tsv	TSV	Spreadsheet editor
Taxa-level.OTU.distribution.combined.png	PNG	Image viewer

6 Formats

Table 18: List of deliverable files, format and recommended programs to access.

Format	Description
TSV	Tab separated table style text file. Can be imported into spreadsheet processing software like MS OFFICE Excel.
PNG	Visual representation in Portable Network Graphics format.

7 Tables

Table 19: Description of filters used.

Name	Thresholds
% Identity	≥ 97.00
E-value	$\leq 1e-06$
% Alignment coverage	≥ 95.00
Min. query length	428
% bitscore threshold for multiple hits	10
Max. hits to consider for multiple hits	50
% abundance	> 0.5

Table 20: Description of Taxonomy.

No.	Name	Description	Example
1	READ.COUNTS	Number of sequences hitting the same OTU	8726
2	PERCENT.COMPOSITION.READS	Percentage of SEQUENCES hitting the same OTU and is calculated based on all the OTUs observed in the sample	21.09
3	CLUSTER.COUNTS	Number of CLUSTERS hitting the same OTU	499
4	PERCENT.COMPOSITION.CLUSTERS	Percentage of CLUSTERS hitting the same OTU and is calculated based on all the OTU CLUSTERS observed in the sample	11.12
5	GI_ID	NCBI GenBank ID	X80725
6	TAXA_ID	NCBI Taxon ID	866789
7	KINGDOM	Name of the kingdom	Bacteria
8	PHYLUM	Name of the phylum	Proteobacteria
9	CLASS	Name of the class	Gammaproteobacteria
10	ORDER	Name of the order	Enterobacteriales
11	FAMILY	Name of the family	Enterobacteriaceae
12	GENUS	Name of the genus	Escherichia
13	SPECIES	Name of the species	Escherichia coli DSM 30083

8 FAQ

Q: What is the necessary coverage for microbiome analysis?

A: The required sequencing depth mainly depends on the complexity of the sample (number, genome size and representation of individual species) and the aim of the project. If you expect your sample to contain only a few different bacteria, a low coverage is sufficient; with many different bacteria expected, a higher coverage is needed. In case of doubt we recommend determining the required depth of sequencing through performing a pilot on a sub-set of samples.

Q: Which organisms can be detected?

A: Phylogenetic characterisation and analysis of microbial communities can be performed for various sample types and organisms. We have tested and demonstrated the utility of this approach for the identification and description of complex and non-complex food/industrial, environmental and medical samples. The focused sequencing of hypervariable regions enables the detection of bacteria present at extremely low frequency.

Q: Down to which taxonomic level can the microbiome be sequenced?

A: Usually the microbiome of a given sample can be resolved down to the genus level with a high degree of certainty. However, related organisms (e.g. belonging to the same genus) may have identical or very similar 16S rRNA genes and therefore, the species cannot be resolved. If the identification of closely related bacteria is of interest, sequencing of further 16S hypervariable regions and/or other genes can be performed.

Q: What is the difference between 'best_hits' and 'multiple_hits'?

A: Both refer to the same BLAST search. While the 'best_hits' summary only takes the first entry of the BLAST hits, the 'multiple_hits' summary utilizes all the hits for the statistics. An additional filter is applied in the 'multiple_hits' evaluation: if a sequence gets more than 50 or 250 hits (depending on the size of the database) these hits are discarded because the informative value is questionable.

Q: How can I open a TSV file in Excel?

A: Start Excel and click File -> Open and select the TSV file you want to open. Next an assistant dialog should show up. Make sure that you select tab as separator. Set the format of all rows without numbers to text. The TSV files use the dot as decimal mark and comma as thousands separator. Make sure that you set both correctly.

Bibliography

- [1] A framework for human microbiome research. *Nature*, 486(7402):215–221, June 2012.
- [2] Tanja Magoč and Steven L. Salzberg. FLASH: Fast Length Adjustment of Short Reads to Improve Genome Assemblies. *Bioinformatics*, 27(21):2957–2963, September 2011.
- [3] W. Li and A. Godzik. Cd-hit: a fast program for clustering and comparing large sets of protein or nucleotide sequences. *Bioinformatics*, 22(13):1658–1659, July 2006.
- [4] Robert C. Edgar, Brian J. Haas, Jose C. Clemente, Christopher Quince, and Rob Knight. UCHIME improves sensitivity and speed of chimera detection. *Bioinformatics*, 27(16):2194–2200, August 2011.
- [5] S. F. Altschul, W. Gish, W. Miller, E. W. Myers, and D. J. Lipman. Basic local alignment search tool. *Journal of molecular biology*, 215(3):403–410, October 1990.
- [6] J. R. Cole, Q. Wang, E. Cardenas, J. Fish, B. Chai, R. J. Farris, A. S. Kulam-Syed-Mohideen, D. M. McGarrell, T. Marsh, G. M. Garrity, and J. M. Tiedje. The Ribosomal Database Project: improved alignments and new tools for rRNA analysis. *Nucleic acids research*, 37(Database issue):D141–D145, January 2009.
- [7] Scott Federhen. The NCBI Taxonomy database. *Nucleic acids research*, 40(Database issue):D136–D143, January 2012.
- [8] Shannon diversity index. http://en.wikipedia.org/wiki/Diversity_index#Shannon_index.
- [9] Simpson diversity index. http://en.wikipedia.org/wiki/Diversity_index#Simpson_index.
- [10] Diversity index. http://en.wikipedia.org/wiki/Diversity_index.
- [11] Ecological Diversity Indices and Rarefaction Species Richness (R package Vegan). <http://cc.oulu.fi/~jarioksa/softhelp/vegan/html/diversity.html>.

GATC Biotech AG
European Genome and Diagnostics Center
Jakob-Stadler-Platz 7
78467 Konstanz

www.gatc-biotech.com