

Microbiome Analysis

31328 Moderne und molekulare Hochdurchsatztechnologien in der medizinischen
Grundlagenforschung

Dr. rer. nat. Israel Barrantes

Research Group Translational Bioinformatics (head)
Institute for Biostatistics and Informatics in Medicine and Ageing Research
Rostock University Medical Center
Ernst-Heydemann-Str. 8, 18057 Rostock, Germany

`israel.barrantes@uni-rostock.de`

October 2025

Overview

- ▶ Setup and introduction to the command line
- ▶ Background on bioinformatics tools
- ▶ Sequence data formats
- ▶ Quality control
- ▶ Sequence assembly
- ▶ Taxonomy assignment

Some Basic Linux Commands

- ▶ **ls** - Lists files and directories in the current location; used to see what's in a folder
- ▶ **cd** - Changes the current directory; used to navigate between folders (e.g., `cd Documents`)
- ▶ **head** - Displays the first few lines of a file; used to preview file contents without opening the entire file
- ▶ **wc** - Counts words, lines, and characters in a file; used to get file statistics (e.g., `wc -l counts lines`)
- ▶ **grep** - Searches for text patterns within files; used to find specific sequences or keywords (e.g., `grep "ATCG" sequences.fasta`)
- ▶ **cat** - Concatenates and displays file contents; used to view entire files or combine multiple files
- ▶ **bash** - The command-line shell/interpreter; used to execute commands and run shell scripts

Why Do We Need the Command Line in Bioinformatics?

- ▶ **Reproducibility:** Command-line tools and programming enable the creation of reproducible workflows. By writing scripts or workflows, researchers can document and automate their analyses, ensuring that others can replicate their results
- ▶ **Pipelines:** Programs talking to each other (pipes)
- ▶ **Redirection:** Programs write and read to files
- ▶ **Text Streams:** Allow us to both couple programs together and process data without storing huge amounts of data in our computers' memory
- ▶ **Modularity:**
 - ▶ Modular workflows allow us to experiment with alternate methods and approaches
 - ▶ Makes it easier to inspect intermediate results and isolate problematic steps
 - ▶ Allows us to choose tools and languages appropriate for specific tasks
 - ▶ Modular programs are reusable and applicable to many types of data

Reference: Buffalo, V. (2015) Bioinformatics Data Skills. O'Reilly Media

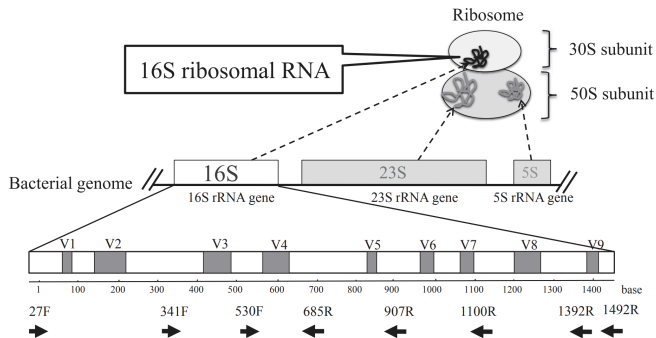
- ▶ **Literate programming:** Chunks of programming (analytical code) with human-readable text (comments)
- ▶ **Version control:** Tracking changes made to sets of files (of a project), typically program source code, scripts and documentation
- ▶ **Environment control:** Versions of all programs (plus libraries, packages, OS) used; archival copies for future reference. Example: `sessionInfo()`
- ▶ **Persistent data sharing:** Collaborative, transparent, accessible science
- ▶ **Documentation:** e.g. README file
- ▶ Project (data + code + ...) validation
- ▶ Command-line tools enable creation of **reproducible workflows** (aka pipelines)
- ▶ **Modularity** of the command-line: Possibility of running long pipelines of programs, one after another + piping

Reference: Ziemann et al. Brief Bioinform. 2023 Sep 22;24(6):bbad375

The 16S rRNA Gene for Community Profiling

- ▶ **Operational concepts of classification:** OTUs and ASVs
 - ▶ **Operational taxonomic units (OTUs):** Consensus sequences from clustering
 - ▶ **Amplicon Sequence Variants (ASVs):** Exact sequences
- ▶ **Community profiling:** Identifying OTUs and/or ASVs in samples
- ▶ 16S rRNA part of the 30S small subunit (SSU) of the prokaryotic ribosome
- ▶ All prokaryotes have one (or more) copy of this gene
- ▶ Different parts of the gene exhibit different levels of conservation
- ▶ Gene ~1540 nt - Not too short to be uninformative; not too long to be unmanageable
- ▶ Most 16S rRNA gene is highly conserved between different species
- ▶ **Hypervariable regions:** Nine much less conserved (V1–V9)

16S rRNA Gene Structure

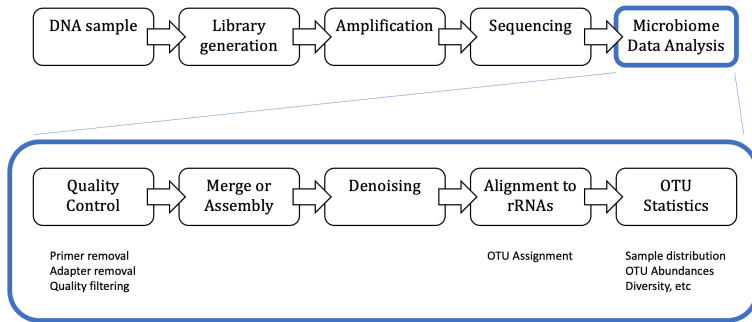


Reference: Fukuda et al. 2016. *J UOEH.*, 38(3):223-32. doi: 10.7888/juoeh.38.223

rRNA Reference Databases

Feature	SILVA	Greengenes	Greengenes2	RDP
Size	~2M seq.	~1.2M seq.	~2.1M seq.	~3M seq.
Last Update	Rel. 138.1 (2023)	13_8 (May 2013)	2022.10 (Oct 2022)	Rel. 11.5 (2016)
Coverage	Bacteria, archaea, euk.	Bacteria, archaea	Bacteria, archaea	Bacteria, archaea
rRNA genes	SSU + LSU	SSU only (16S)	SSU only (16S)	SSU only
Taxonomy	SILVA tax.	Greengenes tax.	GTDB taxonomy	Bergey's Manual
Best For	General use, euk.	Old studies	General use	Old studies

A Typical Microbiome Analysis Experiment



FASTA

```
head -10 egfr_flank.fasta
>ENSMUSG00000020122|ENSMUST00000138518
CCCTCTCATCATGCTGCTCAGTGTATCTCTAAATAGCACTCTCAACCCCGTGAACCTGGT
TATTAAAAACATGCCCAAAGTCTGGGAGCCAGGGCTGCAGGAAATACCAAGCCTCAGT
TCATCAAAACAGTTCATTGCCCAAATGTTCTCAGCTGCAGCTTTCATGAGGTAACCCA
GGGCCCACTGTTCTCTGGT
>ENSMUSG00000020122|ENSMUST00000125984
GAGTCAGGTTGAAGCTGCCCTGAACACTACAGAGAAGAGAGCCCTGGTGCTCTGTTGTC
TCCAGAACCCCAATATGTCTTGTGAAGGGGCACAAACCCCTCAAGGGGTGTCACTTCTT
CTGATCACTTTTGTACTGTTTACTAACTGATCCTATGAATCACTGTGCTTCTCAGAGG
CTGTAAACACACGCTTGCAAT
```

```
>id, description
sequence....
```

FASTQ

```
@DJB775P1:248:D0MDGACXX:7:1202:12362:49613
TGCTTACTCTGCGTTGATACCACTGCTTAGATCGAAGAGCACAGTCTGAA
+
JJJJJIIJJJJJJIHHHHGHHFFFFFEEEEEDBD?DDDDDBDDDBDDCA
@DJB775P1:248:D0MDGACXX:7:1202:12782:49716
CTCTGCGTTGATACCACTGCTTACTCTGCGTTGATACCACTGCTTAGATCGG
+
IIIIIIIIIIIIIIHHHHHHHHFFFFFEECCCBCECCCCCCCCCCCCCCC
```

```
@id, description
sequence
+indicates end of sequence
base quality (Phred format)
```

Viewing the first 10 lines of a FASTQ file:

```
! head Platz1_R1.head.fastq
```

Counting total number of lines:

```
! wc -l Platz1_R1.head.fastq
```

Hint: Each read in FASTQ format consists of four lines

Finding specific nucleotide combinations:

```
! grep "AATATT" Platz1_R1.head.fastq | head
```

Counting occurrences:

```
! grep -c "AATATT" Platz1_R1.head.fastq
```

A quality value Q is an integer mapping of p (i.e., the probability that the corresponding base call is incorrect).

$$Q_{\text{sanger}} = -10 \log_{10} p$$

$$Q_{\text{sanger}} = -10 \log_{10} p$$

$$30 = -10 \log_{10} p$$

$$\frac{30}{-10} = \log_{10} p$$

$$p = 10^{-3}$$

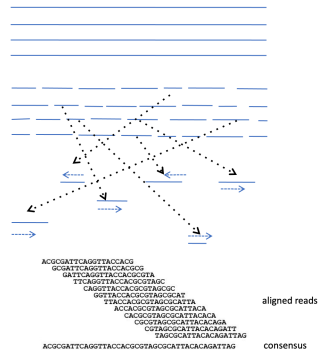
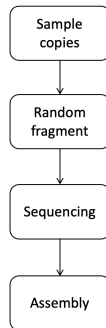
Running FASTQC:

```
! fastqc --quiet Platz1_R1.head.fastq
```

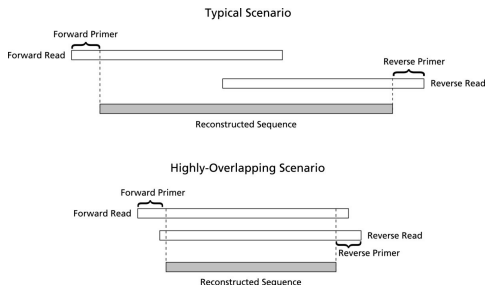
To display the output, download the HTML output from the Files view on the side menu to your local computer (Folder: /course/data2025), and open it in a web browser.

Sequence Assembly

- ▶ Reconstructing complete DNA sequences by aligning and merging smaller overlapping DNA fragments
- ▶ NGS produces short fragments (e.g. 50-300 bp for Illumina)
- ▶ Uses computational algorithms to find overlaps and combine them into longer sequences
- ▶ Analog to solving a puzzle with overlapping pieces



Using PANDAseq for Amplicon Assembly



- ▶ PANDAseq merges forward and reverse reads by finding their overlapping region
- ▶ When reads share identical sequences, it creates a single, longer consensus sequence
- ▶ Uses quality scores to decide which base to keep when there are mismatches
- ▶ Discards reads that do not overlap sufficiently or have too many errors

Reference: Masella et al. *BMC Bioinformatics* 13, 31 (2012). doi:10.1186/1471-2105-13-31

Running PANDAseq:

```
! pandaseq -f Platz10_R1.head.fastq -r Platz10_R2.head.fastq -w Platz10.fa -g  
log.txt
```

Viewing the output:

```
! head Platz10.fa
```

Counting sequences in FASTA output:

```
! grep -c ">" Platz10.fa
```

Question

What is the rate of FASTA sequences vs FASTQ reads? And what does this tell about our sequencing and assembly quality and efficiency?

- ▶ **Community profile:** Represents the taxonomic composition of a microbial community obtained from DNA sequencing data. Typically OTU/ASV tables
- ▶ **Taxonomy Assignment:** Identifying which organism each DNA sequence belongs to by comparing it against reference databases (SILVA, Greengenes, RDP)
- ▶ **OTUs (Operational Taxonomic Units):** Groups of sequences clustered together based on similarity (typically $\geq 97\%$ identity)
- ▶ **ASVs (Amplicon Sequence Variants):** Exact unique sequences identified after error correction, providing single-nucleotide resolution. ASVs are now preferred over OTUs for higher precision and reproducibility

Example OTU Table

samples

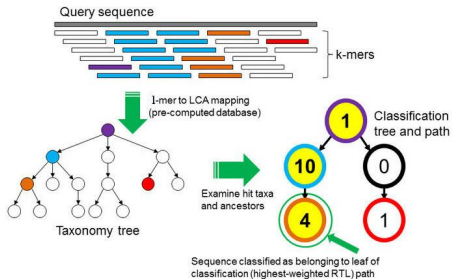
	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U
1	Constructed from bloom file																				
2	OTU ID	Platz10	Platz11	Platz12	Platz13	Platz14	Platz15	Platz16	Platz17	Platz18	Platz19	Platz20	Platz21	Platz22	Platz23	Platz24	Platz25	Platz26	Platz27	taxonomy	
3	10613ee715ee715007f17580e782e1dc	5293	1874	1168	2084	2079	917	2077	815	1843	4594	20519	3310	0	4913	3782	7958	4478	0	d_Bacteria; p_Firmicutes; c_Clostridia; o_Oscillospirales; f_Ruminococcaceae; g_Faecalibacterium	
4	d0a5010d505ce4e2116102789d7520	3066	867	486	1603	2480	0	1370	523	845	2248	8864	2608	0	2821	2093	3935	1777	0	d_Bacteria; p_Firmicutes; c_Clostridia; o_Oscillospirales; f_Ruminococcaceae; g_Faecalibacterium	
5	4376c00a4c90d6a2418a89a40421a57	2930	0	0	839	0	0	0	0	0	0	0	0	0	0	0	0	0	0	d_Bacteria; p_Bacteroidetes; c_Bacteroidia; o_Bacteroidales; f_Bacteroidaceae; g_Bacteroides	
6	8711aaf80a2c8466f4b214a5d43a	2592	0	6708	0	0	185	4	0	0	0	0	0	0	8	1861	0	0	0	d_Bacteria; p_Firmicutes; c_Clostridia; o_Oscillospirales; f_Ruminococcaceae	
7	61227f9b3c16075c50867996a595688	2227	434	1208	0	3566	0	0	0	0	0	0	0	0	1975	0	178	0	0	d_Bacteria; p_Bacteroidetes; c_Bacteroidia; o_Bacteroidales; f_Bacteroidaceae; g_Bacteroides	
8	556ae0ffba4535baab801214519a2e2	1602	853	1332	328	0	45	0	0	0	775	0	136	0	0	11151	0	0	0	d_Bacteria; p_Bacteroidetes; c_Bacteroidia; o_Bacteroidales; f_Bacteroidaceae; g_Bacteroides	
9	9052c040e448425392a69388a71ba83	1529	0	0	5	0	14	4	3	0	0	0	0	14	10	103	32	254	3	20_d_Bacteria; p_Firmicutes; c_Bacilli; o_Lactobacillales; f_Streptococcaceae; g_Lactococcus	
10	58867810629391a679f8335a0993217	1471	0	0	446	0	0	0	0	0	0	0	0	0	0	0	0	0	0	d_Bacteria; p_Bacteroidetes; c_Bacteroidia; o_Bacteroidales; f_Bacteroidaceae; g_Bacteroides	
11	ee151830a70c9ba8a91826a43759a6a2c	1430	211	0	793	295	0	2448	0	49	0	0	0	2460	0	328	2	341	9056	1_d_Bacteria; p_Firmicutes; c_Clostridia; o_Oscillospirales; f_Ruminococcaceae; g_Ruminococcus; s_Rum	
12	0049444c4215c586a7463275c6a4e39	1307	0	4036	0	0	79	2	5	0	0	0	0	0	5	1052	0	0	0	d_Bacteria; p_Firmicutes; c_Clostridia; o_Oscillospirales; f_Ruminococcaceae	
13	c2c37a901b6a788f8c2350a883a1e1c11	1108	257	690	0	4794	0	0	0	0	0	0	0	1663	0	0	0	0	0	d_Bacteria; p_Bacteroidetes; c_Bacteroidia; o_Bacteroidales; f_Bacteroidaceae; g_Bacteroides	
14	10f411747c26cd91eeebbb818c09373	1014	1477	13014	2041	810	1016	595	1742	624	1356	6785	2054	4637	5730	20056	2739	10593	5278	d_Bacteria; p_Firmicutes; c_Clostridia; o_Lachnospirales; f_Lachnospiraceae; g_Agarthobacter	
15	af7c4009d42a5770a97a991c129596a	917	0	0	384	0	0	4762	0	0	0	0	0	0	0	1897	843	195	0	d_Bacteria; p_Firmicutes; c_Clostridia; o_Lachnospirales; f_Lachnospiraceae	
16	3520066a0f7147137a97a2331c1c14f	866	787	1996	239	1082	367	361	1246	2719	371	747	993	742	3337	0	3804	3507	1085	d_Bacteria; p_Firmicutes; c_Clostridia; o_Lachnospirales; f_Lachnospiraceae; g_Blaustia	
17	8a0a867a2670c10ed3f1f466712e5d0d	852	264	439	0	0	7	3	0	2490	4	0	0	0	37	173	2656	109	0	22_d_Bacteria; p_Firmicutes; c_Bacilli; o_Lactobacillales; f_Streptococcaceae; g_Streptococcus	
18	1a0a87a6c0165186e65558b6b931781	806	349	1033	816	1282	3131	940	2417	5160	1325	4523	1264	0	1806	617	1887	2830	1615	d_Bacteria; p_Firmicutes; c_Clostridia; o_Lachnospirales; f_Firmicutes	
19	4dc132748295ae19f1919a2409a9a9	799	0	735	92	0	0	0	0	0	0	0	387	0	191	2	0	6406	0	0_d_Bacteria; p_Bacteroidetes; c_Bacteroidia; o_Bacteroidales; f_Bacteroidaceae; g_Bacteroides	
20	c0d04a4a94b40e4b3775a5a15d077	784	168	2623	285	706	895	496	1863	9051	1777	10190	3131	0	58	12957	3235	9553	19141	d_Bacteria; p_Bacteroidetes; c_Bacteroidia; o_Bacteroidales; f_Bacteroidaceae; g_Bacteroides	
21	0540f712872d71a72150a4c5a131017	775	0	0	3665	0	883	1394	0	788	2085	1360	8054	0	5429	3427	3578	3457	1408	d_Bacteria; p_Firmicutes; c_Clostridia; o_Oscillospirales; f_Ruminococcaceae; g_Subdoligranum	
22	b409725e0cf2131264078648723a9	772	0	3	0	0	0	0	0	0	0	0	0	0	9	0	0	0	0	0_d_Bacteria; p_Firmicutes; c_Bacilli; o_Lactobacillales; f_Streptococcaceae; g_Lactococcus	
23	88c037812707c8a008811eb6a679f	755	0	2164	0	1353	0	0	6477	0	0	0	0	0	297	0	363	0	0	0_d_Bacteria; p_Bacteroidetes; c_Bacteroidia; o_Bacteroidales; f_Bacteroidaceae; g_Bacteroides	
24	0521f48c285174073a20a196869720a	545	4	0	702	0	122	0	1361	4300	5847	0	478	0	2456	0	3606	6131	1385	d_Bacteria; p_Bacteroidetes; c_Bacteroidia; o_Bacteroidales; f_Bacteroidaceae; g_Bacteroides	
25	70104a4e26a51746e04a8e4a5a432	538	0	0	188	0	0	0	3206	0	0	0	0	0	0	0	98	0	0	0_d_Bacteria; p_Firmicutes; c_Clostridia; o_Lachnospirales; f_Lachnospiraceae	
26	4221356a5182004401a193a1030f9e	523	771	7846	1106	962	0	427	1108	3313	0	0	1613	0	3194	11635	0	5147	0	0_d_Bacteria; p_Firmicutes; c_Clostridia; o_Lachnospirales; f_Lachnospiraceae; g_Agarthobacter	
27	b08d24c438024a80c79a1823c3e6d6d	515	0	0	0	0	0	0	0	0	0	0	0	0	0	0	346	0	0	0_d_Bacteria; p_Firmicutes; c_Clostridia; o_Oscillospirales; f_Ruminococcaceae; g_Faecalibacterium	
28	c11402054177a2382c1079206448b	497	188	648	415	1478	1347	563	1399	2060	0	0	958	0	913	0	0	1498	810	d_Bacteria; p_Firmicutes; c_Clostridia; o_Lachnospirales; f_Lachnospiraceae	
29	a5001336a48a630251410051a07d9	469	0	1247	0	1547	0	374	852	1509	155	0	840	3123	3030	0	922	0	0	0_d_Bacteria; p_Firmicutes; c_Clostridia; o_Lachnospirales; f_Lachnospiraceae; g_Blaustia	
30	8033a674a317e0af624a049f740767	451	0	0	2158	0	358	854	0	1045	0	2474	0	2884	1377	0	0	0	0	0_d_Bacteria; p_Firmicutes; c_Clostridia; o_Oscillospirales; f_Ruminococcaceae; g_Subdoligranum	
31	bce774841a641c20f0e9278c7f7f	446	98	1721	201	948	118	373	1363	5446	737	0	2853	0	25	7942	0	5412	10105	d_Bacteria; p_Bacteroidetes; c_Bacteroidia; o_Bacteroidales; f_Bacteroidaceae; g_Bacteroides	
32	107b206972425a2526530707775d0d	438	1562	1516	662	485	388	759	4459	2561	2216	4659	711	8985	6811	1156	5581	1450	4641	d_Bacteria; p_Firmicutes; c_Clostridia; o_Lachnospirales; f_Lachnospiraceae; g_Blaustia	
33	d51740a060f486a4afca11332c23	434	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0_d_Bacteria; p_Bacteroidetes; c_Bacteroidia; o_Bacteroidales; f_Prevotellaceae	
34	1a10f87518760220a617b10a6329f404	406	385	270	0	0	2	2	0	1348	1	0	0	0	136	1311	39	0	0	49_d_Bacteria; p_Firmicutes; c_Bacilli; o_Lactobacillales; f_Streptococcaceae; g_Streptococcus	
35	83a63807d74719a7278939c3a2113c0	403	0	0	245	111	0	395	0	0	0	775	0	0	0	613	362	0	0	0_d_Bacteria; p_Firmicutes; c_Clostridia; o_Oscillospirales; f_Ruminococcaceae; g_Faecalibacterium	
36	c544ab2c6e19743ba64b5c19f10f	384	873	1284	153	949	112	843	227	11224	591	0	0	2	0	798	7088	0	0	0_d_Bacteria; p_Actinobacteriota; c_Actinobacteria; o_Bifidobacteriales; f_Bifidobacteriaceae; g_Alistipes	
37	d00a3a6829e441b0d1f720551518f	344	349	0	860	502	449	153	150	4162	885	0	2266	0	0	0	6	1150	0	0_d_Bacteria; p_Bacteroidetes; c_Bacteroidia; o_Bacteroidales; f_Tannerellaceae; g_Parabacteroides; s_Pa	
38	4c43703a613203292c31281b0a2c5	342	0	0	402	0	44	0	8062	2298	3062	0	335	0	1366	0	1885	3105	0	0_d_Bacteria; p_Bacteroidetes; c_Bacteroidia; o_Bacteroidales; f_Bacteroidaceae; g_Bacteroides	
39	35793272a0a11288a40f79322a	333	0	1514	0	1842	0	0	4639	0	0	0	246	0	273	0	0	0	0	0_d_Bacteria; p_Bacteroidetes; c_Bacteroidia; o_Bacteroidales; f_Bacteroidaceae; g_Bacteroides	
40	467676dc764587187b119171a73a4942a	301	0	0	76	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0_d_Bacteria; p_Bacteroidetes; c_Bacteroidia; o_Bacteroidales; f_Bacteroidaceae; g_Bacteroides	
41	040891317053a4c1151955f1076c438	293	33	0	501	488	518	340	51	0471	7498	1091	1044	0	0	0	0	136	1861	0_d_Bacteria; p_Bacteroidetes; c_Bacteroidia; o_Bacteroidales; f_Bifidobacteriaceae; g_Alistipes	

taxa

Taxonomy Assignment with Kraken2

Kraken2 uses the k-mer method:

- ▶ Reference database with sequences cut into k-mers (e.g., 30 bases)
- ▶ Break new sequence into k-mers
- ▶ Check each k-mer in database
- ▶ Assign to lowest common ancestor (LCA) if found in multiple organisms
- ▶ Decide organism based on where most k-mers point



Reference: Wood & Salzberg, 2014. *Genome Biol* 15: R46. doi:10.1186/gb-2014-15-3-r46

Running Kraken2 against Greengenes database:

```
! kraken2 --db 16S_Greengenes_k2db --use-names --output output.txt --report  
  report.txt --paired Platz10_R1.head.fastq Platz10_R2.head.fastq
```

Inspecting results:

```
! cat report.txt
```

Question

What are the most predominant genera in your personal Illumina runs?

The output report from Kraken2 consists of the following fields:

- ▶ Percentage of fragments covered by the clade rooted at this taxon
- ▶ Number of fragments covered by the clade rooted at this taxon
- ▶ Number of fragments assigned directly to this taxon
- ▶ A rank code: (U)nclassified, (R)oot, (D)omain, (K)ingdom, (P)hylum, (C)lass, (O)rder, (F)amily, (G)enus, or (S)pecies
- ▶ NCBI taxonomic ID number
- ▶ Indented scientific name

Execute Kraken2 with all Illumina sequencing pairs:

```
! bash getkrakenreports.sh
```

Merge all reports with kraken-biom:

```
! kraken-biom *.report --fmt tsv -o mbtmicrobiome20251022.tsv
```

View the OTU table:

```
! cat mbtmicrobiome20251022.tsv
```

OTU taxonomic information can be retrieved from the NCBI Taxonomy database

Aims

- ▶ How do the samples group with each other? Principal Components
- ▶ How much diversity is there in a sample? Alpha diversity
- ▶ How does the diversity differ between samples? Beta diversity
- ▶ What is the taxonomic content of the sample?
- ▶ Rarefaction: Adjusting for differences in library sizes across samples to aid comparisons

Reference: Willis A. 2019. Front. Microbiol. 10: 2407, doi:10.3389/fmicb.2019.02407

Why R for Bioinformatics?

- ▶ Freely available
- ▶ User-friendly resources to learn R
- ▶ R editors and integrated development environments (IDEs)
- ▶ Active user community
- ▶ Large code repositories: CRAN, Bioconductor, github, gitlab
- ▶ Ready-made packages and functions for Bioinformatics
- ▶ Reproducibility: Code, library versions, Rmarkdown, Jupyter notebooks

Install Bioconductor and phyloseq:

```
install.packages("BiocManager")  
BiocManager::install("phyloseq")
```

Load required libraries:

```
# load ggplot2 library (graphics)  
library(ggplot2)  
  
# loading phyloseq library (microbiome analysis)  
library(phyloseq)
```

Loading Input Files

```
# OTU data
InputBiomFile <- "mbtmicrobiome2024.biom"

# Samples' metadata
InputMapFile <- "sample-metadata-2024.tsv"

# prepare phyloseq object by loading both files
BiomData <- import_biom(InputBiomFile, parseFunction = parse_taxonomy_
  greengenes)
SampleData <- import_qiime_sample_data(InputMapFile)
```

View metadata:

```
head(SampleData)
```

Creating a Phyloseq Object

Merge OTU and sample data:

```
# create phyloseq object by merging OTU and sample data  
ExperimentPhyloseqObject <- merge_phyloseq(BiomData, SampleData)
```

Check object properties:

```
# checking the features of our original microbiome data  
ExperimentPhyloseqObject
```

Subsetting and Filtering Data

```
# create a temporary phyloseq object for working
psTemp <- ExperimentPhyloseqObject

# Prune OTUs with low abundances from all samples
psTemp <- prune_taxa(taxa_sums(psTemp) > 100, psTemp)

# Prune samples with no metadata
psTemp <- subset_samples(psTemp, Gender != "U")

# checking the features of our microbiome data
psTemp
```

Question

What are the differences between the ExperimentPhyloseqObject and psTemp objects?

- ▶ Ordination is a visualization technique that places samples in 2D or 3D space based on their microbial community similarity
- ▶ Samples with similar microbiomes appear close together; different ones are far apart
- ▶ Ordination shows quickly if samples cluster by treatment, disease state, or other factors—revealing biological patterns in your data at a glance
- ▶ **PCoA (Principal Coordinates Analysis):** Works with ecological distances like Bray-Curtis dissimilarity, which better captures compositional differences between communities

References: Bray & Curtis (1957). Ecol Monogr, doi:10.2307/1942268

Ramette (2007). FEMS Microbiol Ecol, doi:10.1111/j.1574-6941.2007.00375.x

Hands-on: Sample Ordination

```
# Calculate distance and ordination
iDist <- distance(psTemp, method="bray")
iMDS  <- ordinate(psTemp, distance=iDist)

# plot sample ordination, e.g. by Gender
plot_ordination(psTemp, iMDS, color="Gender")
```

Question

Are there any clear separations between the gender groups? Why/Why not?

Ordination by other metadata (e.g. "Pet"):

```
plot_ordination(psTemp, iMDS, color="Pet") +
  geom_text(aes(label=SampleID), vjust = -1)
```

Phylum and Genus Barplots

- ▶ Microbial communities are usually represented by bar plots showing the relative abundance (percentage) of different bacterial groups
- ▶ Each color represents a different taxon, and bar height shows how much of the community it comprises
- ▶ After taxonomy assignment, sequences are grouped by taxonomic level and counted. Counts are converted to percentages so all bars add up to 100%
- ▶ Provides an intuitive visual summary of community composition
- ▶ Typically only the top 10-20 most abundant taxa are shown individually; rare taxa are grouped as "Other"

References: Caporaso et al. (2010). Nature Methods, doi:10.1038/nmeth.f.303

McMurdie & Holmes (2013). PLoS One, doi:10.1371/journal.pone.0061217

- ▶ Normalization technique that randomly subsamples all samples to the same sequencing depth (number of reads) to enable fair comparison
- ▶ Different samples often have different total read counts due to technical variation in sequencing
- ▶ All samples are randomly downsampled to match the sample with the fewest reads
- ▶ **However:** Some argue rarefaction discards valuable data and recommend alternative normalization methods (e.g., CSS, DESeq2), particularly for differential abundance testing

References: Hughes et al. (2001). Appl Environ Microbiol, doi:10.1128/AEM.67.10.4399-4406.2001
McMurdie & Holmes (2014). PLoS Comput Biol, doi:10.1371/journal.pcbi.1003531

Absolute abundances at the Phylum and Genus levels:

```
### Phylum  
plot_bar(psTemp, "SampleID", fill="Phylum")  
  
### Genus  
plot_bar(psTemp, "SampleID", fill="Genus")
```

Note

These plots show that the samples are not directly comparable. Hence rarefaction is needed to adjust for differences in library sizes.

Hands-on: Abundances After Rarefaction

```
# initializes the random number generator to ensure reproducibility
set.seed(1212)

# Rarefaction to an even depth
ps.rarefied <- rarefy_even_depth(psTemp)

# Agglomerate the data to a single taxonomic level, e.g. Phylum
ps.rarefied.glom <- tax_glom(ps.rarefied, "Phylum")

# Plot abundances
plot_bar(ps.rarefied.glom, "SampleID", fill="Phylum")
```

Question

What are the differences between the abundance plots before and after rarefaction?

```
### Abundances per category, e.g. "Gender"
# Merge samples by the selected category
mergedGP <- merge_samples(psTemp, "Gender")

# Rarefaction to an even depth
ps.rarefied <- rarefy_even_depth(mergedGP)

# Agglomerate to Phylum level
ps.rarefied.glom <- tax_glom(ps.rarefied, "Phylum")

# Plot abundances for the example category "Gender"
plot_bar(ps.rarefied.glom, fill="Phylum")
```

- ▶ **Diversity:** How many different species are present and how evenly distributed they are in microbial communities
- ▶ **Alpha diversity:** Diversity **within** a single sample
 - ▶ Common metrics: species richness, Shannon index, observed ASVs/OTUs
 - ▶ Shannon index combines richness and evenness: higher values = more diverse communities
- ▶ **Beta diversity:** Diversity **between** samples
 - ▶ Measured using distance metrics like Bray-Curtis dissimilarity or UniFrac distances
 - ▶ Used to create ordination plots (PCoA, PCA)
- ▶ Low alpha diversity may indicate dysbiosis; high beta diversity between groups suggests different conditions are shaping distinct communities

References: Magurran (2004). Measuring Biological Diversity, doi:10.1002/9780470999738
Lozupone & Knight (2005). Appl Environ Microbiol, doi:10.1128/AEM.71.12.8228-8235.2005

Shannon diversity on individual samples:

```
plot_richness(psTemp, x = "SampleID", measures = c("Shannon"))
```

Shannon diversity by Gender:

```
plot_richness(psTemp, x = "Gender", color = "Gender", measures = c("Shannon"))
```

Enhanced Diversity Visualization

```
# Assign our plot to a variable
Our_Richness_plot <- plot_richness(psTemp, x = "Gender",
                                   color = "Gender",
                                   measures = c("Shannon"))

# Improving our plot by adding features
Our_Richness_plot +
  geom_boxplot(data = Our_Richness_plot$data,
              aes(x = Gender, y = value, color = Gender),
              alpha = 0.1) + # boxplot
  labs(title = "Richness of the gut microbiome (Shannon alpha diversity)",
       subtitle = "MBT Class 2024") + # title and subtitle
  theme(axis.text.x = element_text(angle = 0, hjust = 0.5)) # x-axis labels
```

Quiz: Steps for the Analysis of Microbiome Data

Statement	T	F
Alignment to rRNA sequences from standard databases	T	F
microRNA identification	T	F
RNA structure prediction	T	F
Quality control	T	F
Denoising or quality filtering	T	F

Quiz: Assembly in Microbiome Analysis Means?

Statement	T	F
Discard sequences with poor FASTQ quality	T	F
Comparison (alignment) of the amplicon sequences versus known rRNAs from ribosomal rRNA databases	T	F
Preparing OTU tables by counting the number of rRNA matches present for each taxa	T	F
Merging fragments (in FASTQ format) into whole sequences (in FASTA format)	T	F
Building a small database of ribosomal RNAs	T	F

Quiz: Abundance (Composition) Plots Show:

Statement	T	F
The phyla and genera compositions simultaneously for all samples in the same plot	T	F
The genera (genus) and species composition simultaneously of a microbial community	T	F
The abundance of all taxonomic levels at the same time for every individual sample	T	F
The abundance of all taxonomic levels at the same time for every sample group	T	F
The composition of a given single taxonomic level for a microbial community	T	F

Quiz: R for Bioinformatics

Statement	T	F
Requires a commercial licence for use	T	F
Large code repositories	T	F
Jupyter Notebooks are limited to Python and hence cannot be used with the R programming language	T	F
Ready-made packages and functions for bioinformatics	T	F
Reproducibility depends on R library versions	T	F

Quiz: Microbial Communities

Statement	T	F
Rarefaction is a method that adjusts for differences in library sizes across samples to aid comparisons	T	F
The principal components method is used to visualize how samples group together	T	F
Genera composition tells the abundance of species in a sample or group of samples	T	F
The taxonomic content of a sample (or group) can only be studied using known taxa	T	F
Alpha diversity measures can be seen as a summary statistic of a single population (within-sample diversity)	T	F