

Microbiome workshop – data generation and processing

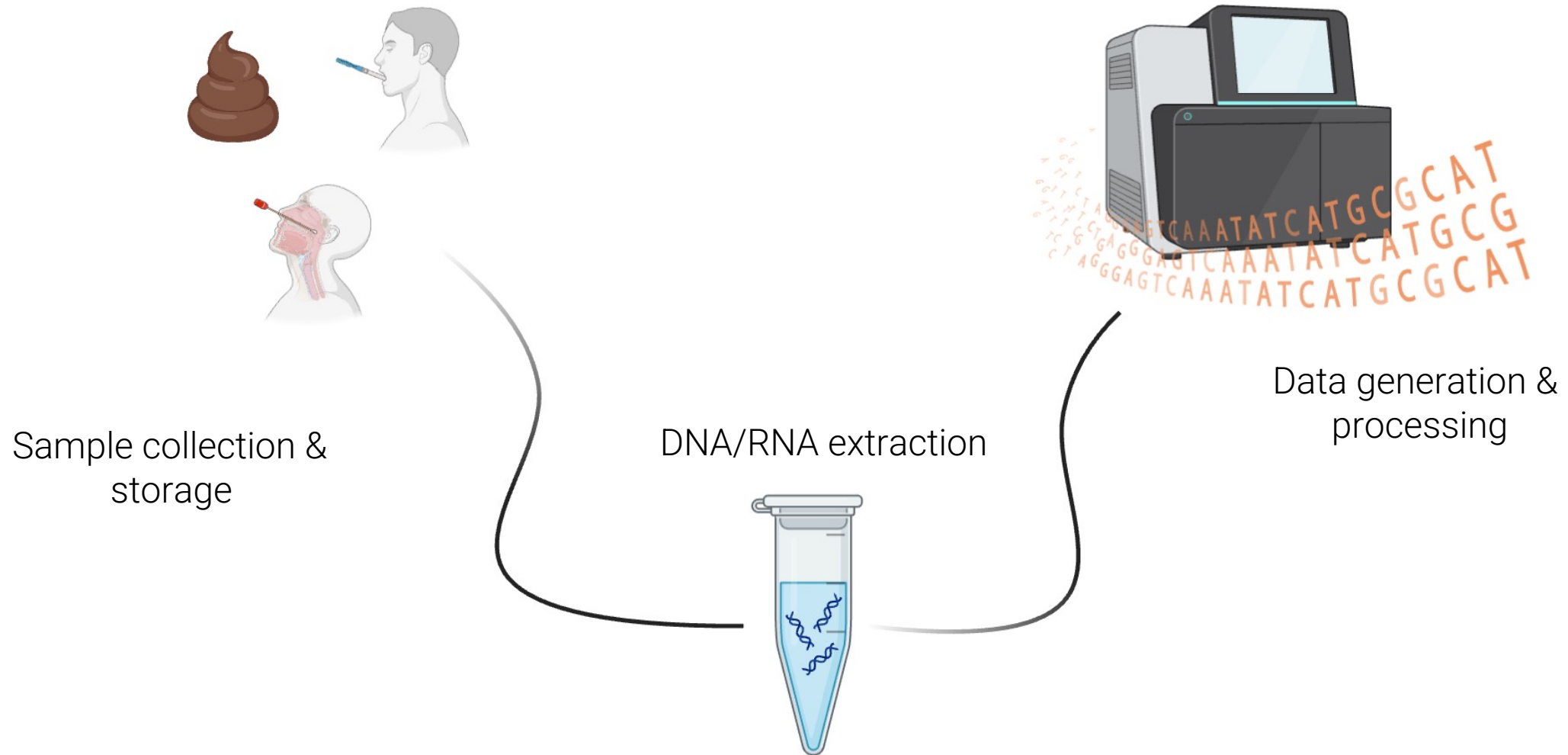
September 2024



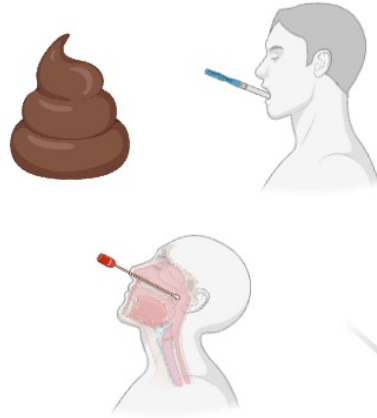
IKMB

Institute of Clinical
Molecular Biology

Microbiome science - from sample to data



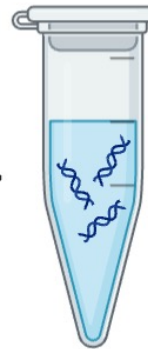
Microbiome science - from sample to data



Sample collection & storage

- Sampling method (swab, tube, ...)
- Sampling buffer (RNAlater, glycerol, ...)
- Storage conditions (snap freezing, -20/-80°C)

DNA/RNA extraction



- Extraction kit used
- DNA or RNA or both?
- Cell disruption (Proteinases, Bead beating)

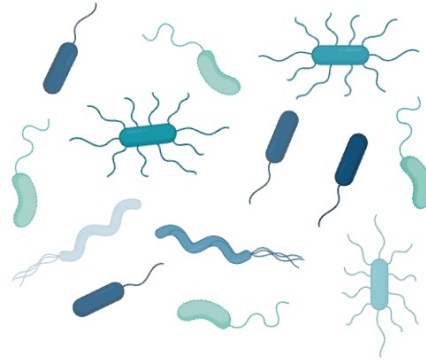


Data generation & processing

- 16S rRNA gene or Shotgun metagenomics
- Short reads or long reads
- Data annotation

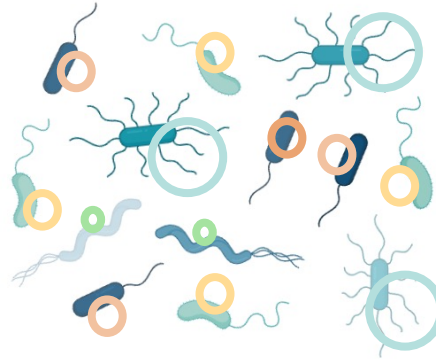
16S rRNA gene
amplicon sequencing

shotgun
metagenomics



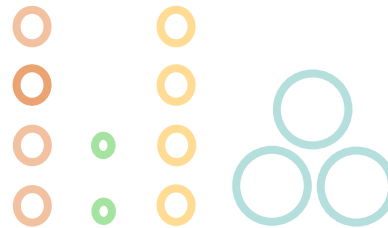
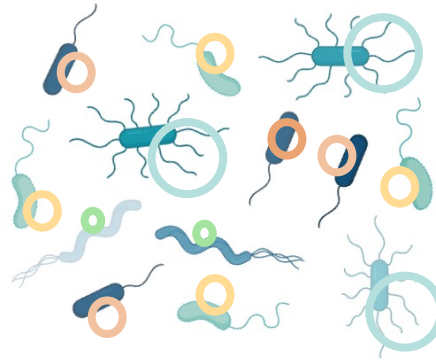
16S rRNA gene
amplicon sequencing

shotgun
metagenomics



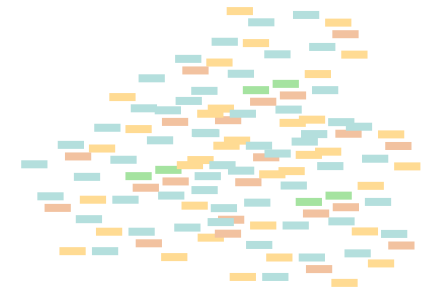
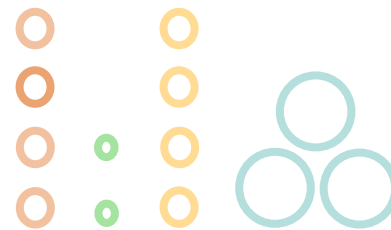
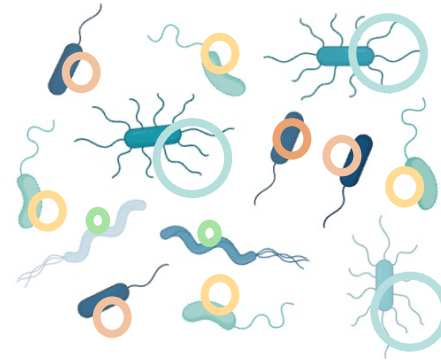
16S rRNA gene
amplicon sequencing

shotgun
metagenomics



16S rRNA gene amplicon sequencing

shotgun metagenomics



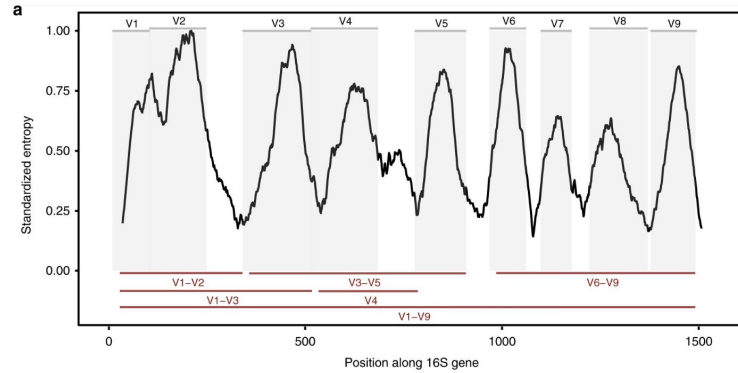
Full genome information + + +

More data + / -

More expensive - -

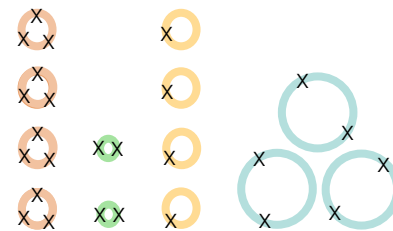
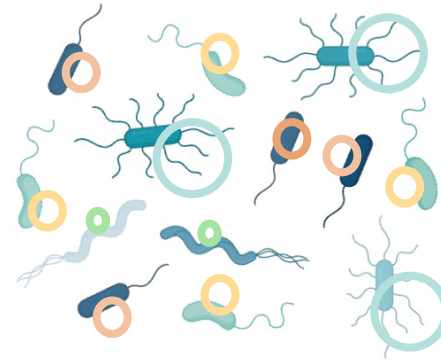
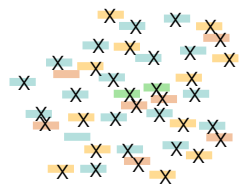
Bias towards high abundant taxa -

16S rRNA gene amplicon sequencing



16S rRNA gene: conserved and variable regions

PCR target + Molecular clock
= taxon information



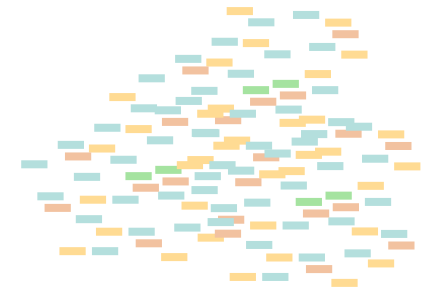
shotgun metagenomics

Full genome information + + +

More data + / -

More expensive - -

Bias towards high abundant taxa -



What we have: The FastQ format

Identifier — | @HWI-EAS209_0006_FC706VJ:5:58:5894:21141#ATCACG/1
Sequence — | TTAATTGGTAAATAAATCTCCTAATAGCTTAGATNTTACCTTNNNNNNNNNNNTAGTTTCTTGAGA
+ sign & identifier — | +HWI-EAS209_0006_FC706VJ:5:58:5894:21141#ATCACG/1
Quality scores — | efcfffffcfeefffcfffffddf`feed]`_]_Ba_^__[YBBBBBBBBBBRTT\]][] dddd`

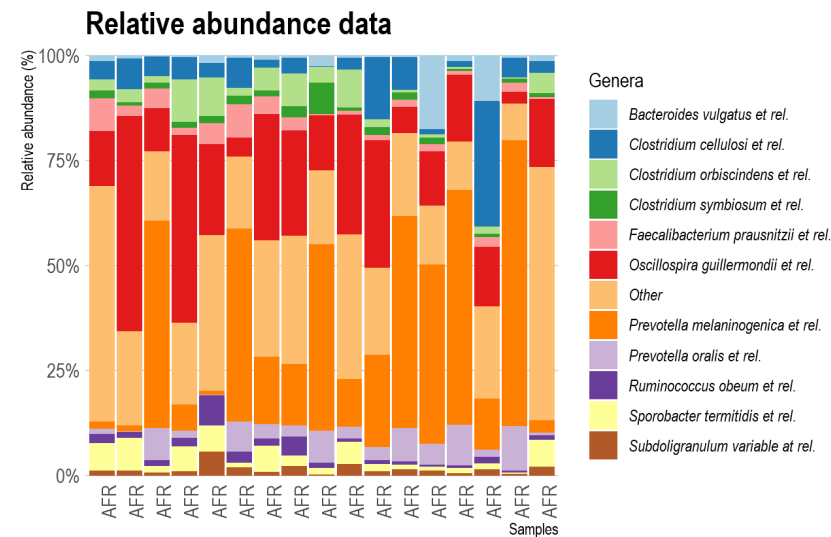
Base T
phred Quality] = 29

What we have: The FastQ format

Identifier — | @HWI-EAS209_0006_FC706VJ:5:58:5894:21141#ATCACG/1
Sequence — | TTAATTGGTAAATAAATCTCCTAATAGCTTAGATNTTACCTTNNNNNNNNNNNTAGTTTCTTGAGA
+ sign & identifier — | +HWI-EAS209_0006_FC706VJ:5:58:5894:21141#ATCACG/1
Quality scores — | efcfffffcfeefffcfffffddf`feed]`_]_Ba_^__[YBBBBBBBBBBRTT\]][] dddd`

Base T
phred Quality] = 29

What we want:



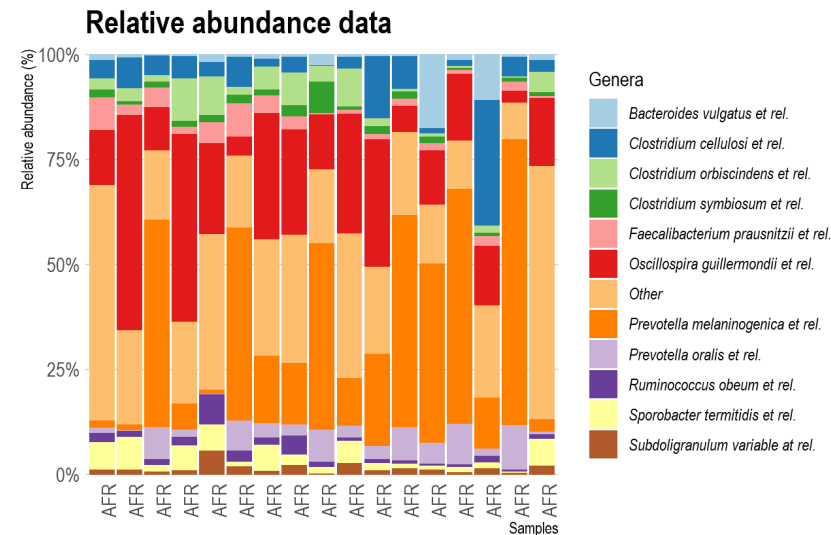
<https://compgenomr.github.io/book/fasta-and-fastq-formats.html>

What we have: The FastQ format

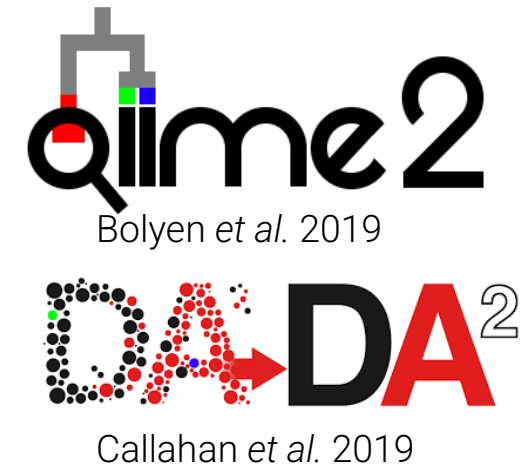
Identifier — | @HWI-EAS209_0006_FC706VJ:5:58:5894:21141#ATCACG/1
Sequence — | TTAATTGGTAAATAAATCTCCTAATAGCTTAGATNTTACCTTNNNNNNNNNNNTAGTTTCTTGAGA
+ sign & identifier — | +HWI-EAS209_0006_FC706VJ:5:58:5894:21141#ATCACG/1
Quality scores — | efcfffffcfeefffcfffffdff`feed]`_]_Ba_^__[YBBBBBBBBBBRTT\]][] dddd`

Base T
phred Quality] = 29

What we want:

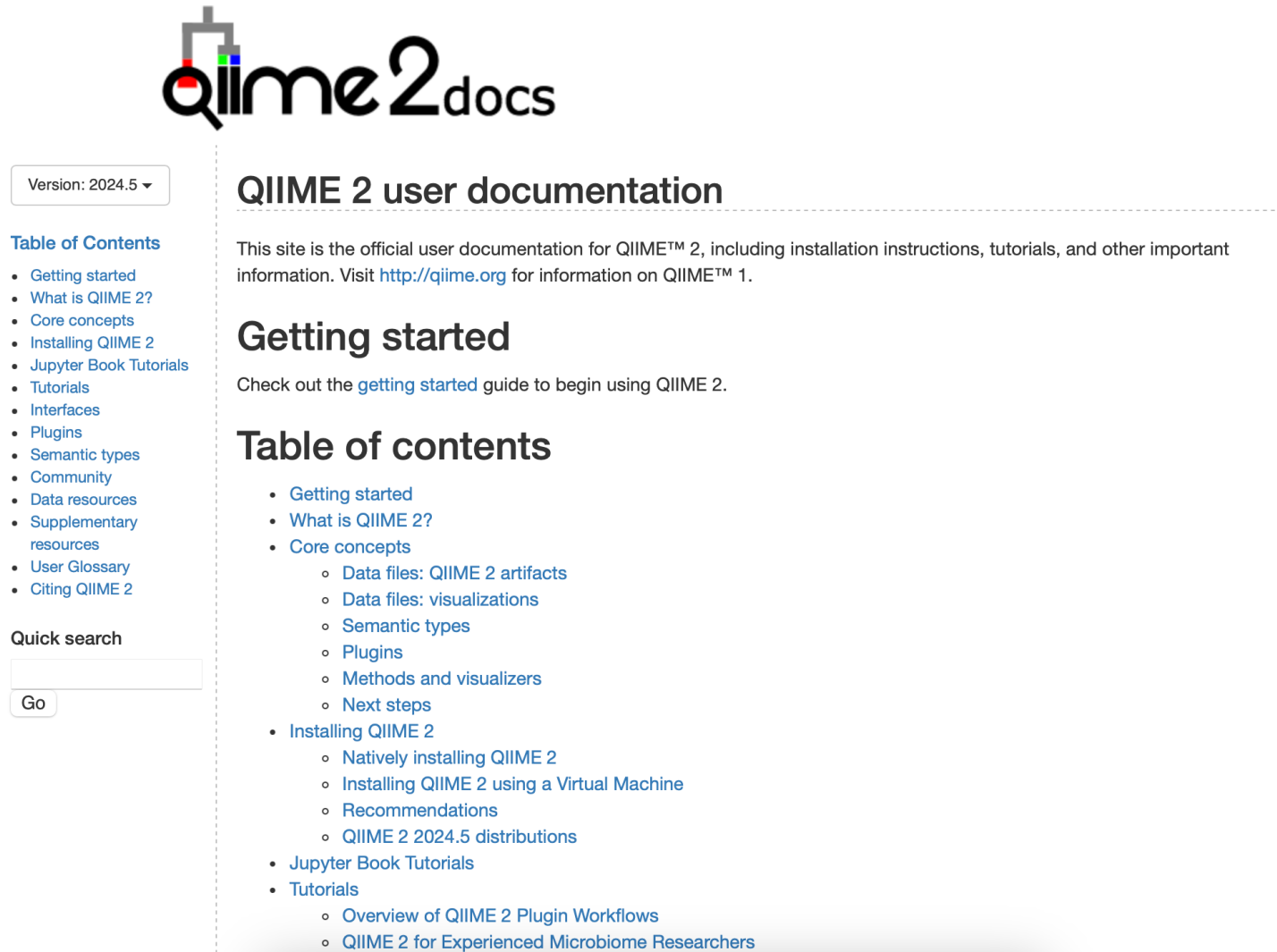


How we get there:



<https://compgenomr.github.io/book/fasta-and-fastq-formats.html>

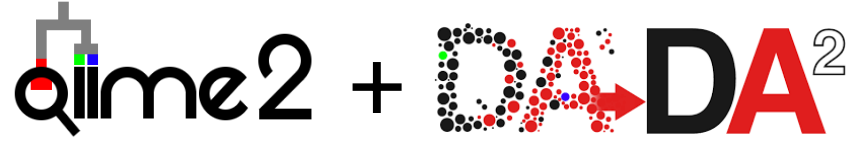
16S data workflow with



The screenshot shows the QIIME 2 user documentation website. At the top, the logo "qiime2docs" is displayed, with "qiime" in a stylized font and "2docs" in a standard font. Below the logo, a dropdown menu shows "Version: 2024.5". The main content area is titled "QIIME 2 user documentation" and contains a paragraph stating that the site is the official user documentation for QIIME™ 2, including installation instructions, tutorials, and other important information. It also mentions that users can visit <http://qiime.org> for information on QIIME™ 1.

On the left side, there is a "Table of Contents" section with a list of links: "Getting started", "What is QIIME 2?", "Core concepts", "Installing QIIME 2", "Jupyter Book Tutorials", "Tutorials", "Interfaces", "Plugins", "Semantic types", "Community", "Data resources", "Supplementary resources", "User Glossary", and "Citing QIIME 2". Below this is a "Quick search" section with a text input field and a "Go" button.

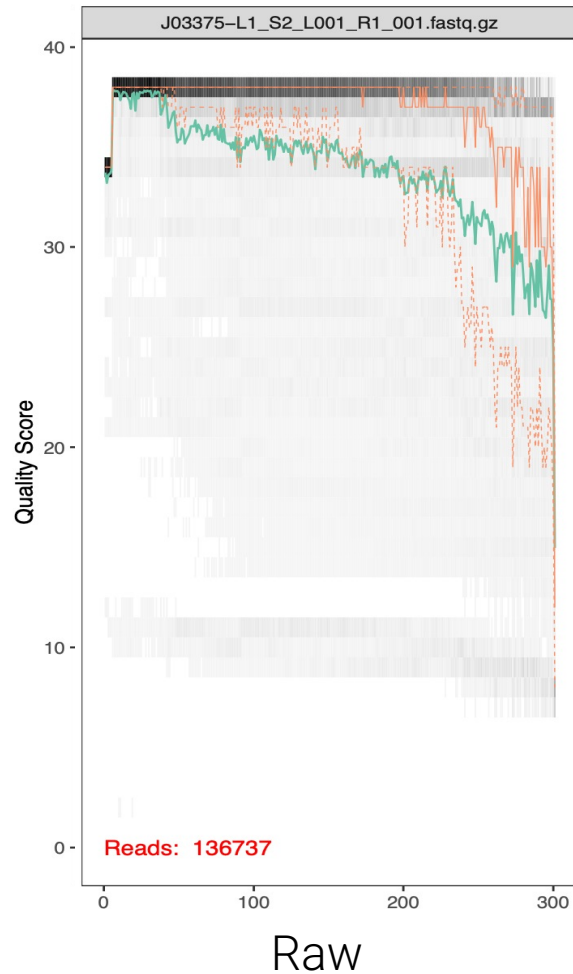
The main content area also features a "Getting started" section with a link to the "getting started" guide to begin using QIIME 2. Below this is a "Table of contents" section with a list of links: "Getting started", "What is QIIME 2?", "Core concepts", "Installing QIIME 2", "Jupyter Book Tutorials", and "Tutorials". The "Core concepts" and "Installing QIIME 2" sections have sub-links: "Data files: QIIME 2 artifacts", "Data files: visualizations", "Semantic types", "Plugins", "Methods and visualizers", "Next steps", "Natively installing QIIME 2", "Installing QIIME 2 using a Virtual Machine", "Recommendations", and "QIIME 2 2024.5 distributions". The "Tutorials" section has sub-links: "Overview of QIIME 2 Plugin Workflows" and "QIIME 2 for Experienced Microbiome Researchers".



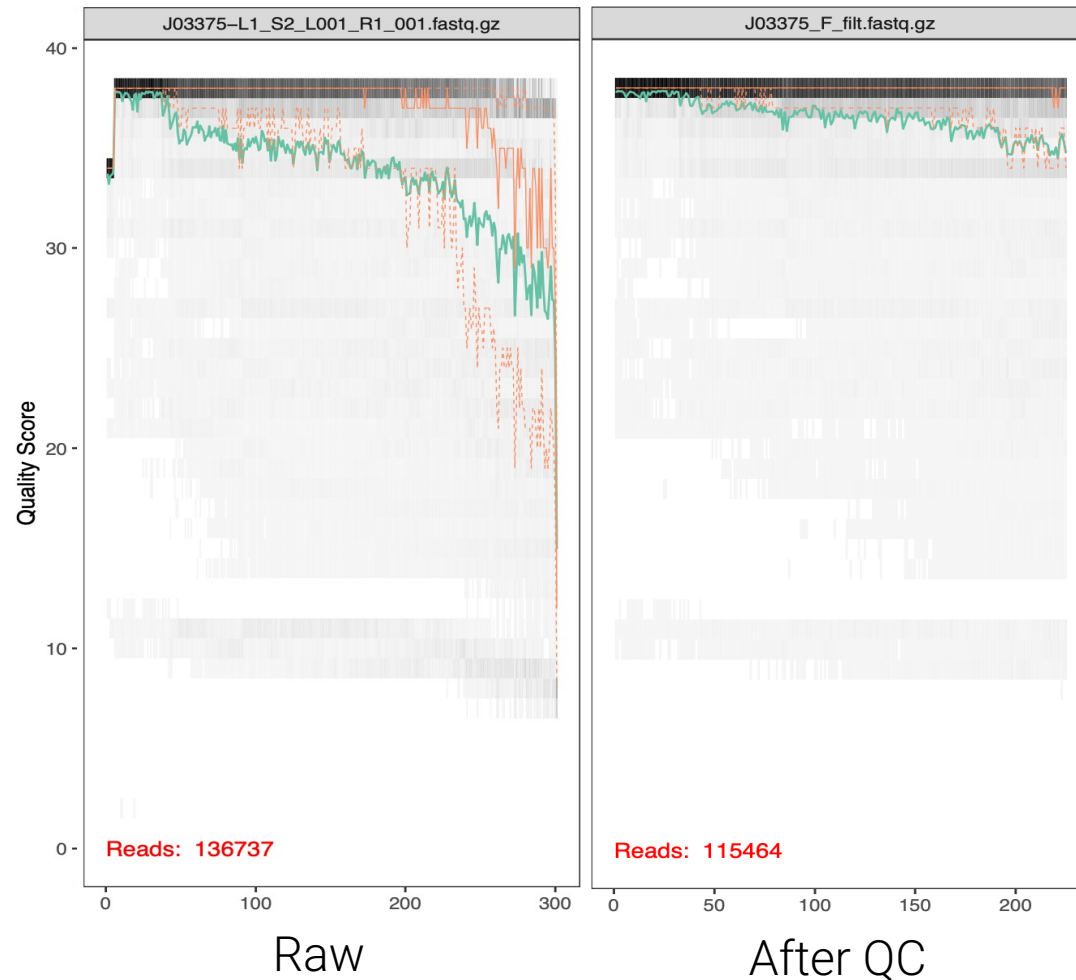
- Tutorials

- Overview of QIIME 2 Plugin Workflows
- QIIME 2 for Experienced Microbiome Researchers
- "Moving Pictures" tutorial
- "Moving Pictures" tutorial - Multiple Interface Edition
- Fecal microbiota transplant (FMT) study: an exercise
- "Atacama soil microbiome" tutorial
- Parkinson's Mouse Tutorial
- Importing data
- Exporting data
- Metadata in QIIME 2
- Filtering data
- Training feature classifiers with q2-feature-classifier
- Evaluating and controlling data quality with q2-quality-control
- Predicting sample metadata values with q2-sample-classifier
- Performing longitudinal and paired sample comparisons with q2-longitudinal
- Identifying and filtering chimeric feature sequences with q2-vsearch
- Alternative methods of read-joining in QIIME 2
- Clustering sequences into OTUs using q2-vsearch
- Utilities in QIIME 2
- Phylogenetic inference with q2-phylogeny

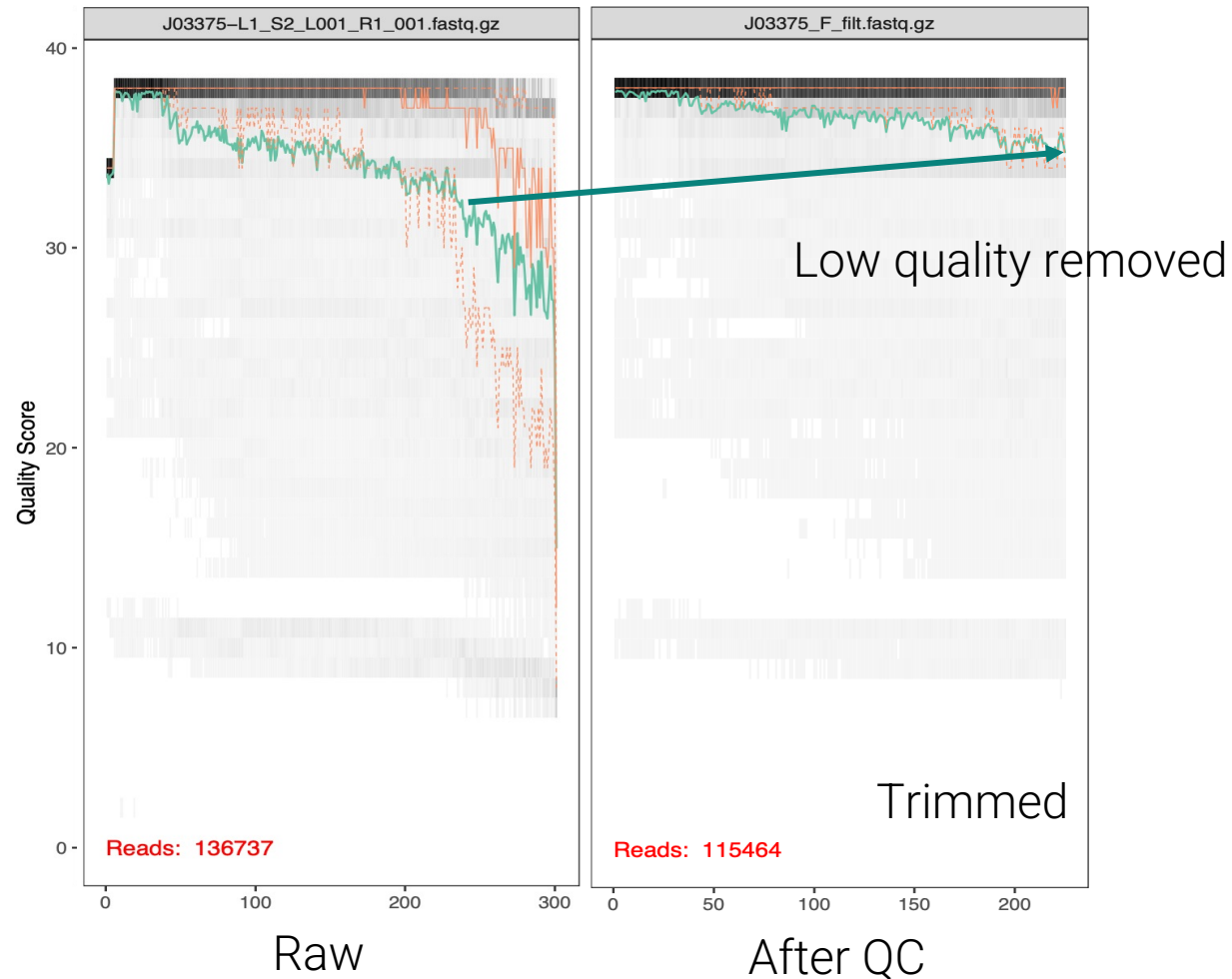
Step 1: QC and trimming



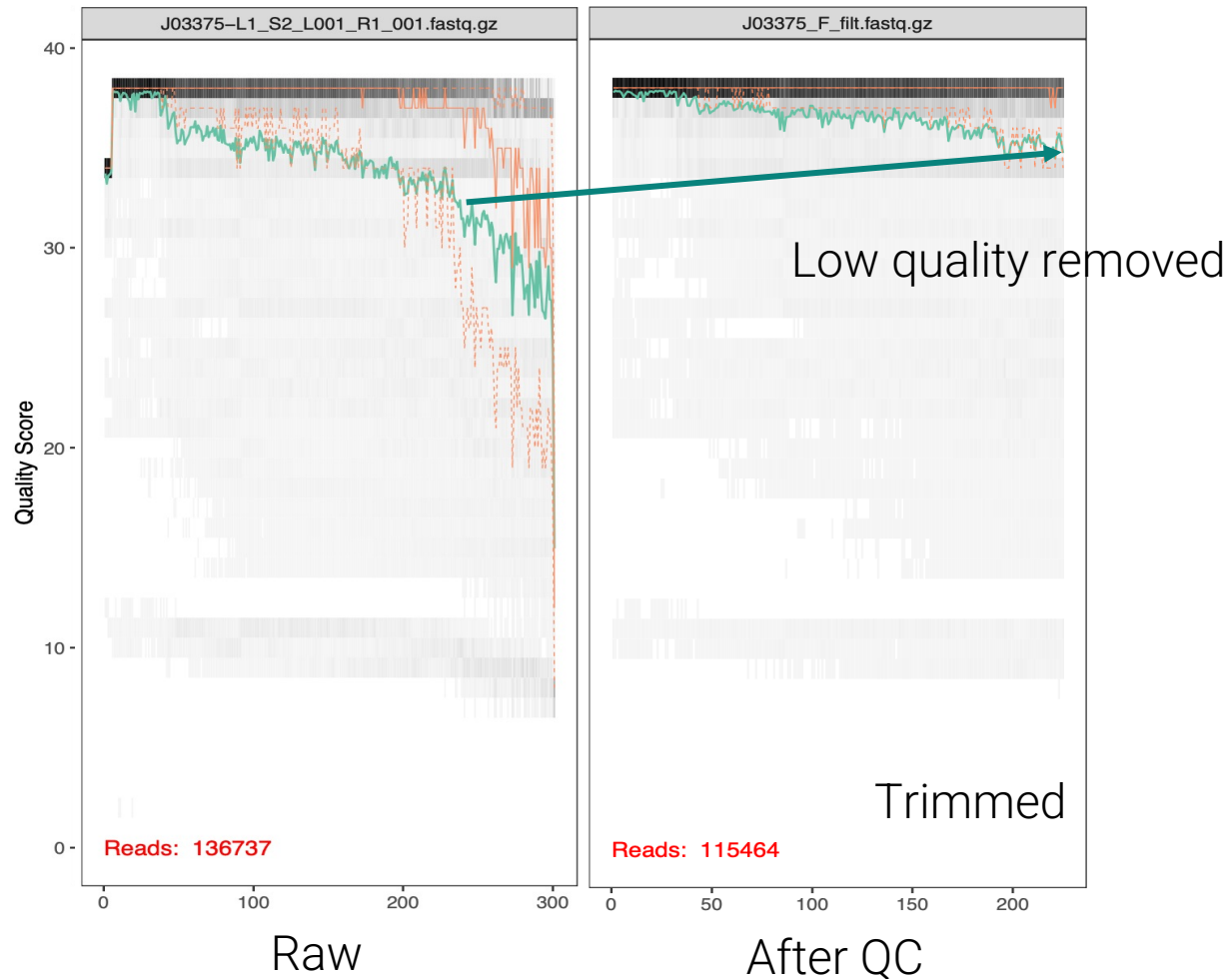
Step 1: QC and trimming



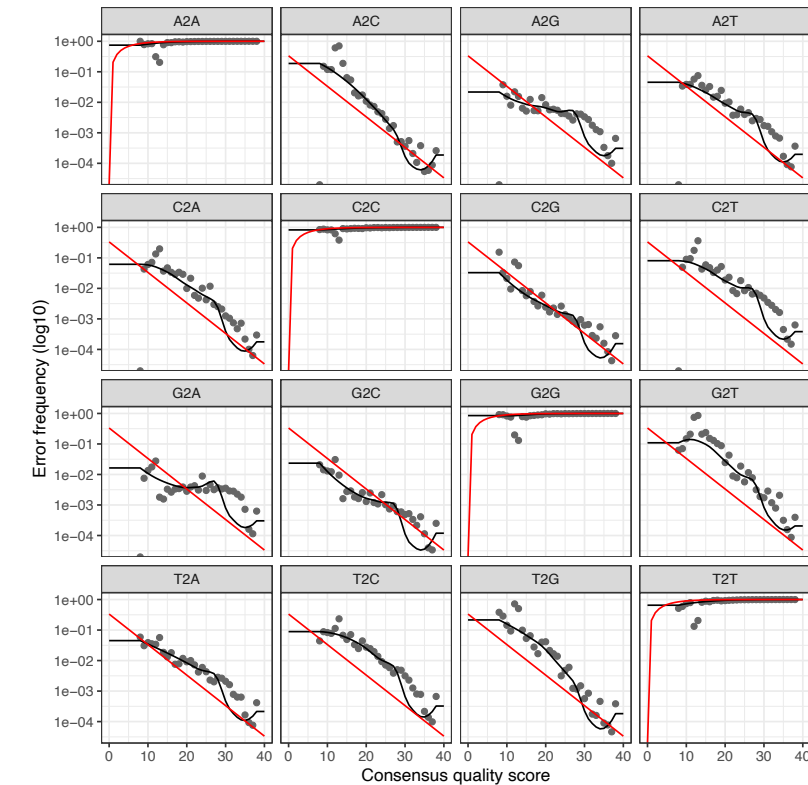
Step 1: QC and trimming



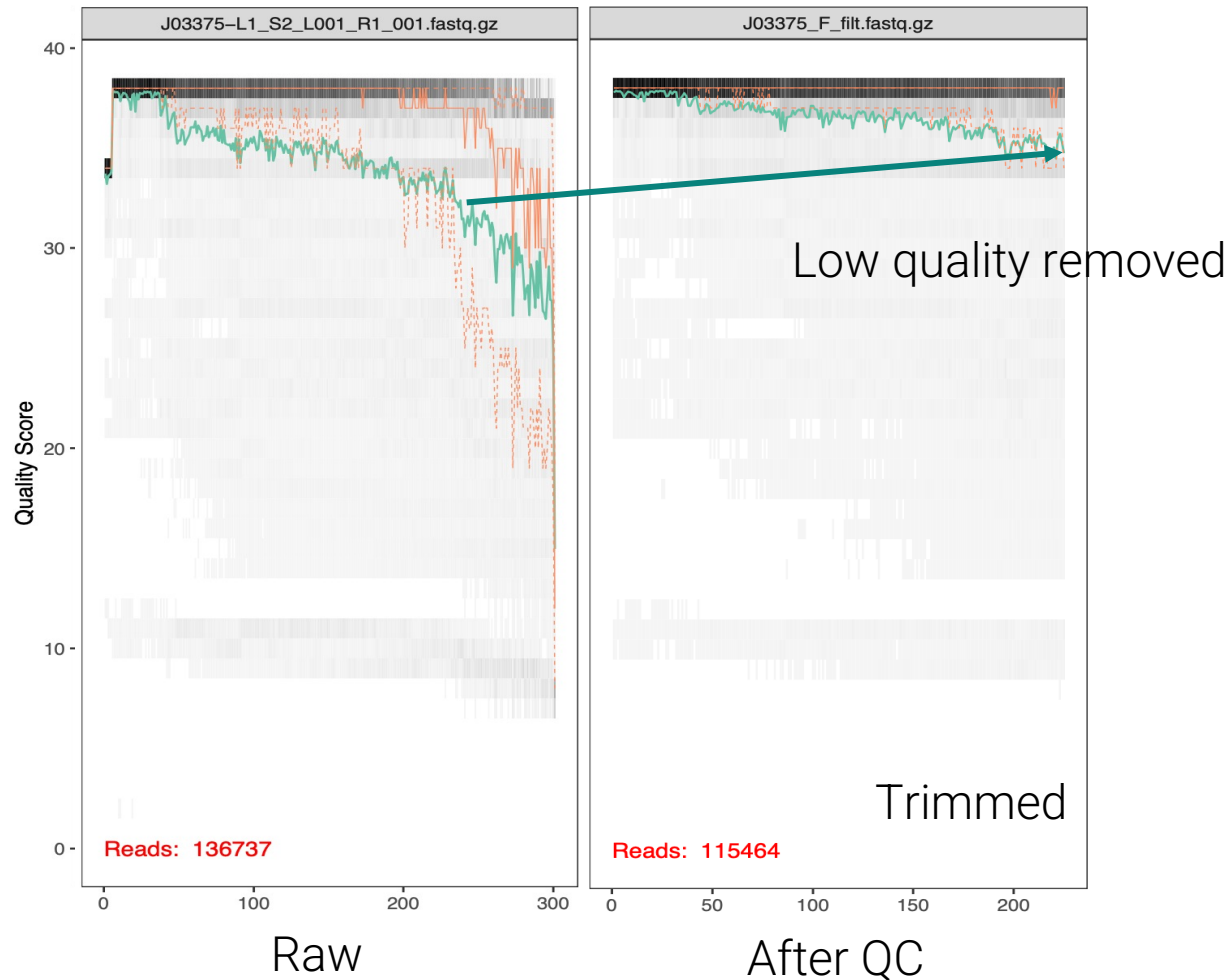
Step 1: QC and trimming



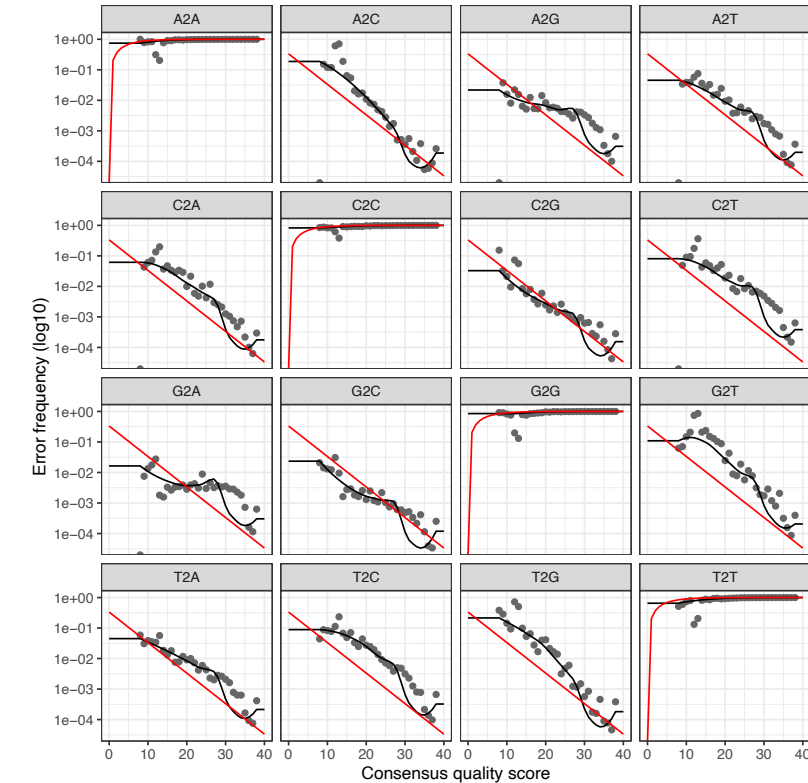
Step 2: Error correction



Step 1: QC and trimming



Step 2: Error correction

Step 3:
Inference of Amplicon Sequence Variants (ASVs)

Assigning taxonomic labels

Abundance Table

	ASV1	ASV2	ASV3	ASV4
S1	5	10	3	7
S2	5	0	3	0
S3	0	0	4	3

ASV sequences:

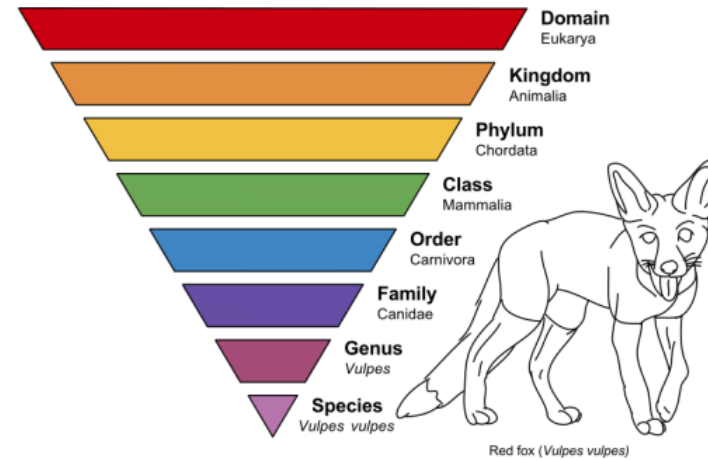
```
>ASV1
ACGCTGGCGGCATGCTTAACACATGCAAGTCGTACGAATGAAT...
>ASV2
ACGCTGGCGGTATGCCTAACACATGCAAGTCGAACGAGGTAGC...
>ASV3
ACGCTAGCGGCAGGCTTAACACATGCAAGTCGAGGGGTAACAG...
>ASV4
ACGCTAGCGGCAGGCTTAACACATGCAAGTCGAGGGGTAGTCA...
```

Assigning taxonomic labels

What is a taxonomy?

Abundance Table

	ASV1	ASV2	ASV3	ASV4
S1	5	10	3	7
S2	5	0	3	0
S3	0	0	4	3



ASV sequences:

```
>ASV1
ACGCTGGCGGCATGCTTAACACATGCAAGTCGTACGAATGAAT...
>ASV2
ACGCTGGCGGTATGCCTAACACATGCAAGTCGAACGAGGTAGC...
>ASV3
ACGCTAGCGGCAGGCTTAACACATGCAAGTCGAGGGGTAAACAG...
>ASV4
ACGCTAGCGGCAGGCTTAACACATGCAAGTCGAGGGGTAGTCA...
```

Assigning taxonomic labels

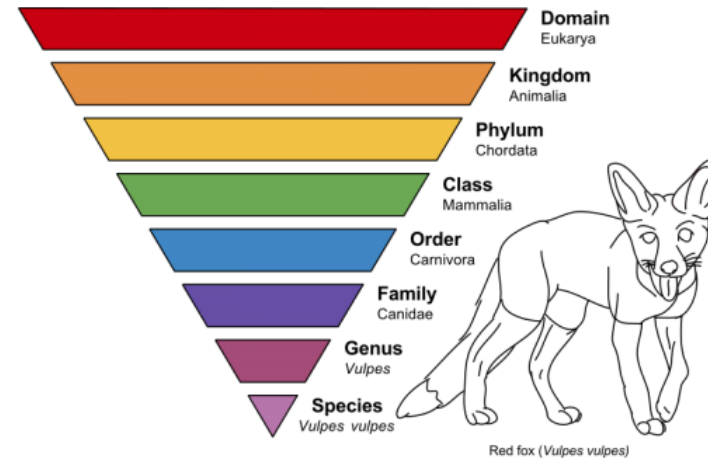
Abundance Table

	ASV1	ASV2	ASV3	ASV4
S1	5	10	3	7
S2	5	0	3	0
S3	0	0	4	3

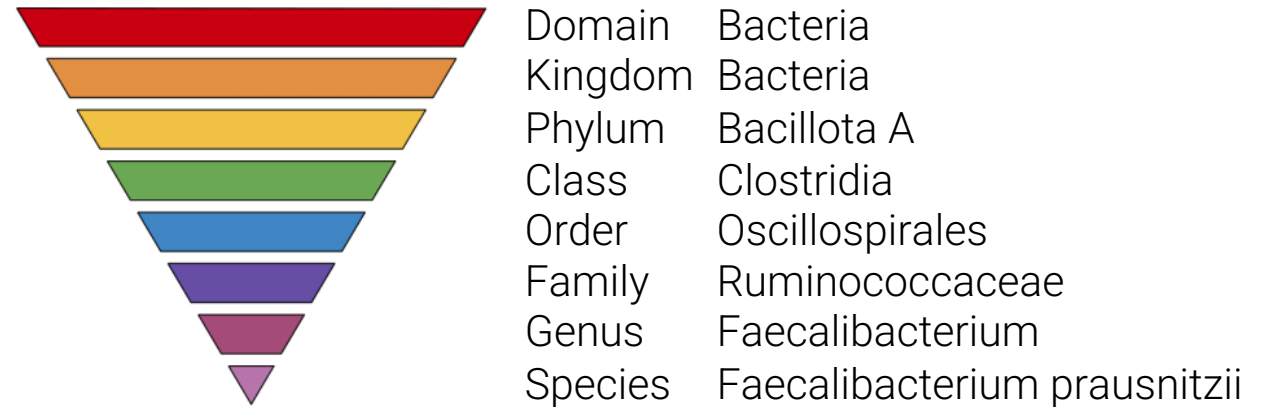
ASV sequences:

```
>ASV1
ACGCTGGCGGCATGCTTAACACATGCAAGTCGTACGAATGAAT...
>ASV2
ACGCTGGCGGTATGCCTAACACATGCAAGTCGAACGAGGTAGC...
>ASV3
ACGCTAGCGGCAGGCTTAACACATGCAAGTCGAGGGGTAAACAG...
>ASV4
ACGCTAGCGGCAGGCTTAACACATGCAAGTCGAGGGGTAGTCA...
```

What is a taxonomy?



Same can be done for microorganisms:



Assigning taxonomic labels

Abundance Table

	ASV1	ASV2	ASV3	ASV4
S1	5	10	3	7
S2	5	0	3	0
S3	0	0	4	3

ASV sequences:

```
>ASV1
ACGCTGGCGGCATGCTTAACACATGCAAGTCGTACGAATGAAT...
>ASV2
ACGCTGGCGGTATGCCTAACACATGCAAGTCGAACGAGGTAGC...
>ASV3
ACGCTAGCGGCAGGCTTAACACATGCAAGTCGAGGGGTAACAG...
>ASV4
ACGCTAGCGGCAGGCTTAACACATGCAAGTCGAGGGGTAGTCA...
```

Database search
(Naïve Bayesian Classifier)



Assigning taxonomic labels

Abundance Table

	ASV1	ASV2	ASV3	ASV4
S1	5	10	3	7
S2	5	0	3	0
S3	0	0	4	3

ASV sequences:

```
>ASV1
ACGCTGGCGGCATGCTTAACACATGCAAGTCGTACGAATGAAT...
>ASV2
ACGCTGGCGGTATGCCTAACACATGCAAGTCGAACGAGGTAGC...
>ASV3
ACGCTAGCGGCAGGCTTAACACATGCAAGTCGAGGGGTAACAG...
>ASV4
ACGCTAGCGGCAGGCTTAACACATGCAAGTCGAGGGGTAGTCA...
```



Ribosomal
database project



Database search
(Naïve Bayesian Classifier)

Assigning taxonomic labels

Abundance Table

	ASV1	ASV2	ASV3	ASV4
S1	5	10	3	7
S2	5	0	3	0
S3	0	0	4	3

ASV sequences:

```
>ASV1
ACGCTGGCGGCATGCTTAACACATGCAAGTCGTACGAATGAAT...
>ASV2
ACGCTGGCGGTATGCCTAACACATGCAAGTCGAACGAGGTAGC...
>ASV3
ACGCTAGCGGCAGGCTTAACACATGCAAGTCGAGGGGTAACAG...
>ASV4
ACGCTAGCGGCAGGCTTAACACATGCAAGTCGAGGGGTAGTCA...
```



Greengenes2



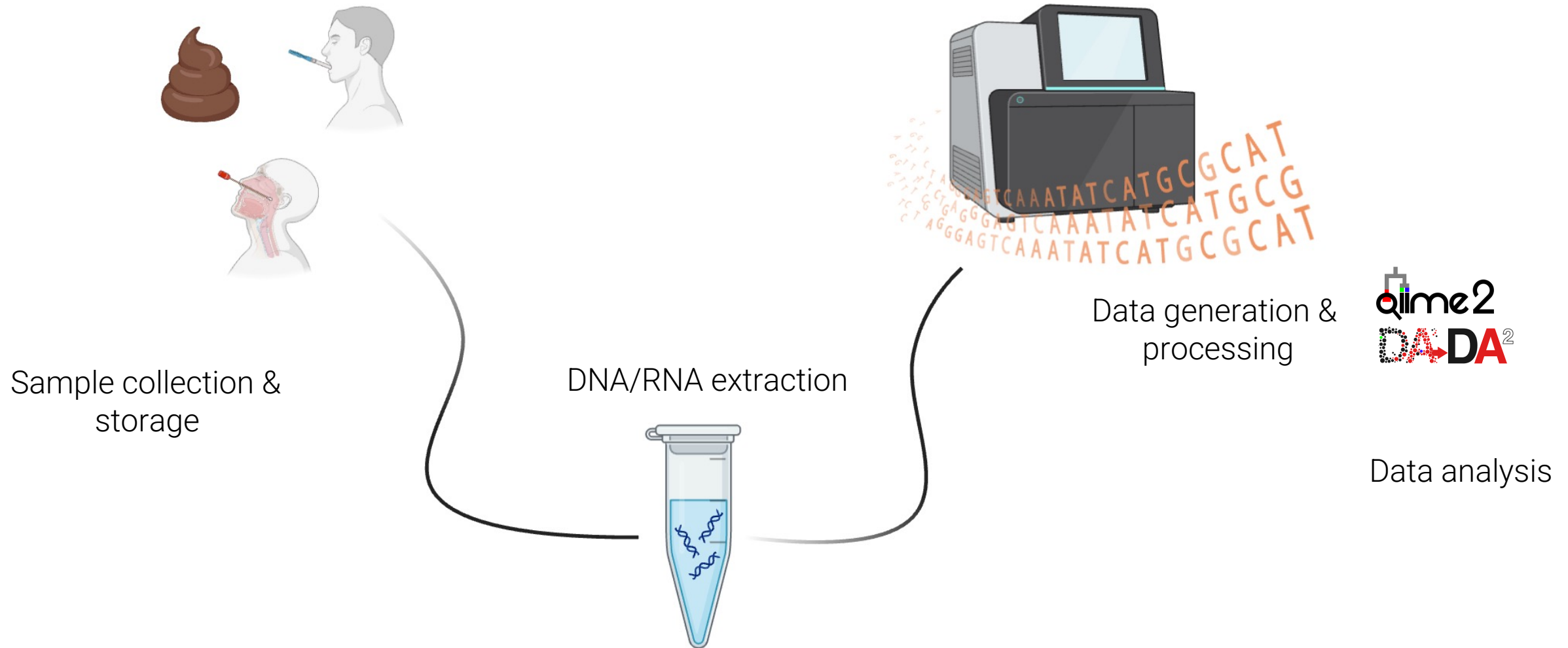
Ribosomal
database project



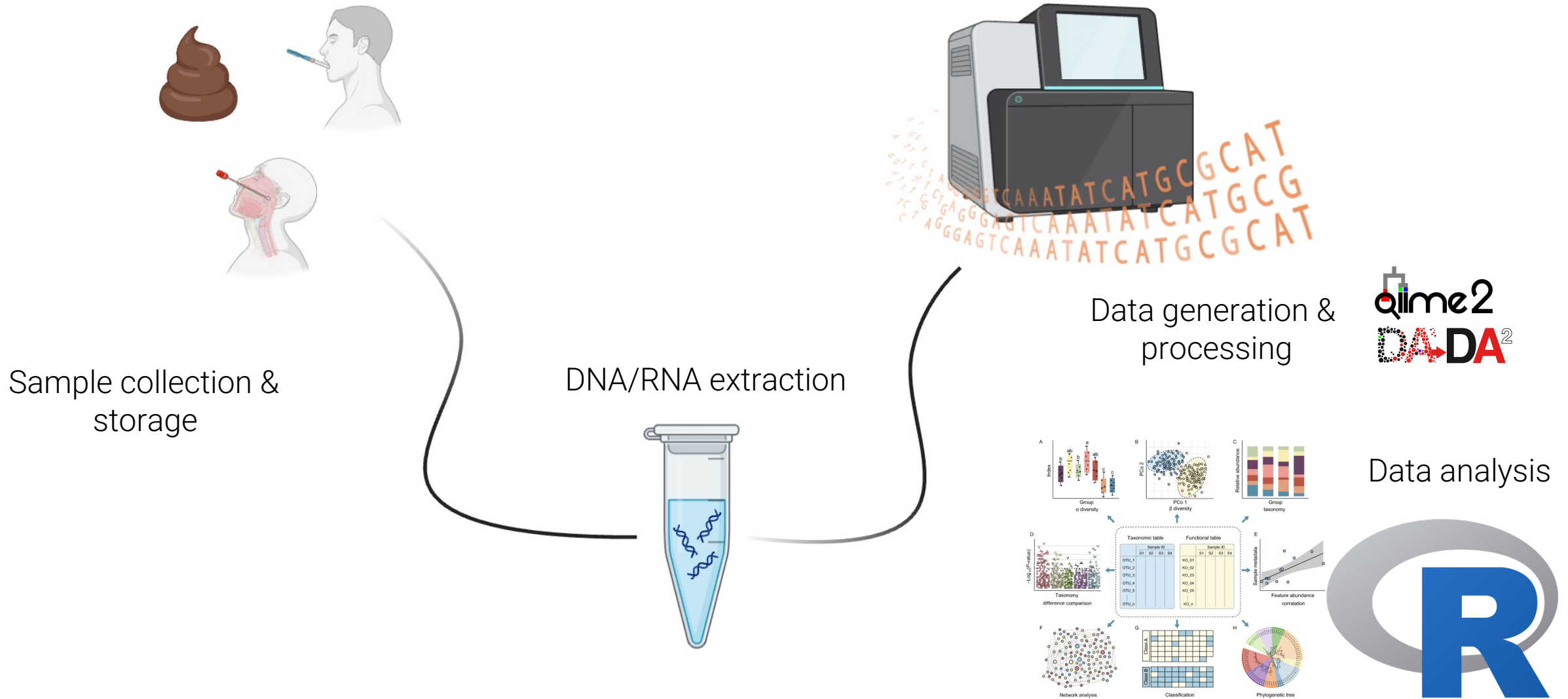
Database search
(Naïve Bayesian Classifier)

	Domain	Phylum	Class	Order	Family	Genus
ASV1	Bacteria	Bacillota_C	Negativicutes	Veillonellales	Dialisteraceae	Dialister
ASV2	Bacteria	Bacteroidota	Bacteroidia	Bacteroidales	Bacteroidaceae	Bacteroides
ASV3	Bacteria	Pseudomonadota	Gammaproteo bacteria	Enterobacterales	Enterbacteriaceae	
ASV4	Bacteria	Bacteroidota	Bacteroidia	Bacteroidales	Bacteroidaceae	Prevotella

Microbiome science - from sample to data



Microbiome science - from sample to data to analysis



RStudio

Go to file/function | Addins | Project: (None)

AlphaDiversity.R | Untitled3* | Intro_R.R | 1_IgA_CoatingIndexCalculation_OT. >>

Source on Save | Run | Source

```

1
2  ## My microbiome analyses
3
4  library(phyloseq)
5
6  setwd("/Users/apschaan/Documents/Kiel/UKSH-CAU/Classes/KMC_workshop")
7
8  taxTable <- read.csv("taxonomy_gmbc1_form.csv", header=TRUE, row.names=1)
9
10 head(taxTable)
11
12
13
14

```

Environment | History | Connections

Import Dataset | 994 MiB | List | Taxa

Data

- TaxaForHeatm... 40 obs. of 12 variables
- TaxaForHeatm... 20 obs. of 2 variables

Values

taxa	taxa_names	taxaKeep
"a361d6b053ce9d38182128a316635e3e"	chr [1:197] "2e63e8e3f9f0f102f570e7..."	chr [1:2174] "9492bc05eee3b253c1557..."

Files | Plots | Packages | Help | Viewer | Presentation

Zoom | Export

Wilcoxon, $p = 4.3e-05$

log(|C|)

Industrialized Non-industrialized

Console | Terminal | Background Jobs

R 4.3.1 ~ /Documents/Kiel/UKSH-CAU/Classes/KMC_workshop/

```

> library(phyloseq)
> setwd("/Users/apschaan/Documents/Kiel/UKSH-CAU/Classes/KMC_workshop")
> taxTable <- read.csv("taxonomy_gmbc1_form.csv", header=TRUE, row.names=1)
> head(taxTable)

```

	Kingdom	Phylum
675c847bccbc53942ebb7b8cbb4efc4d	d__Bacteria	p__Bacteroidota
da6a7cba87e0895b9cbe6037b9bd8b3b	d__Bacteria	p__Bacteroidota