

IMPACTT-MIC



From Sequences to ASV Tables

Rob Beiko and Diana Haider

(with thanks to Jacob Nearing and others)

Dalhousie University





Learning Objectives

- By the end of this lecture, you will be able to:
 - **Understand** the main strengths and weaknesses of the 16S rRNA gene
 - **Read and interpret** the contents of sequence files
 - **Describe** the process of sequence clustering
 - **Understand** the key differences between OTUs, ASVs, and taxonomic summaries



The Plan

- (1) Why 16S? (10 minutes)
- (2) Sequence data and quality control (10 minutes)
- (3) Clustering your sequences (25-ish minutes)

Why 16S?



Taxonomy

- The classification of organisms according to a predefined hierarchy
- “How do we assign individuals” has often been based on what was measurable at the time (morphology, biochemistry, DNA hybridization, marker-gene similarity, whole-genome similarity)



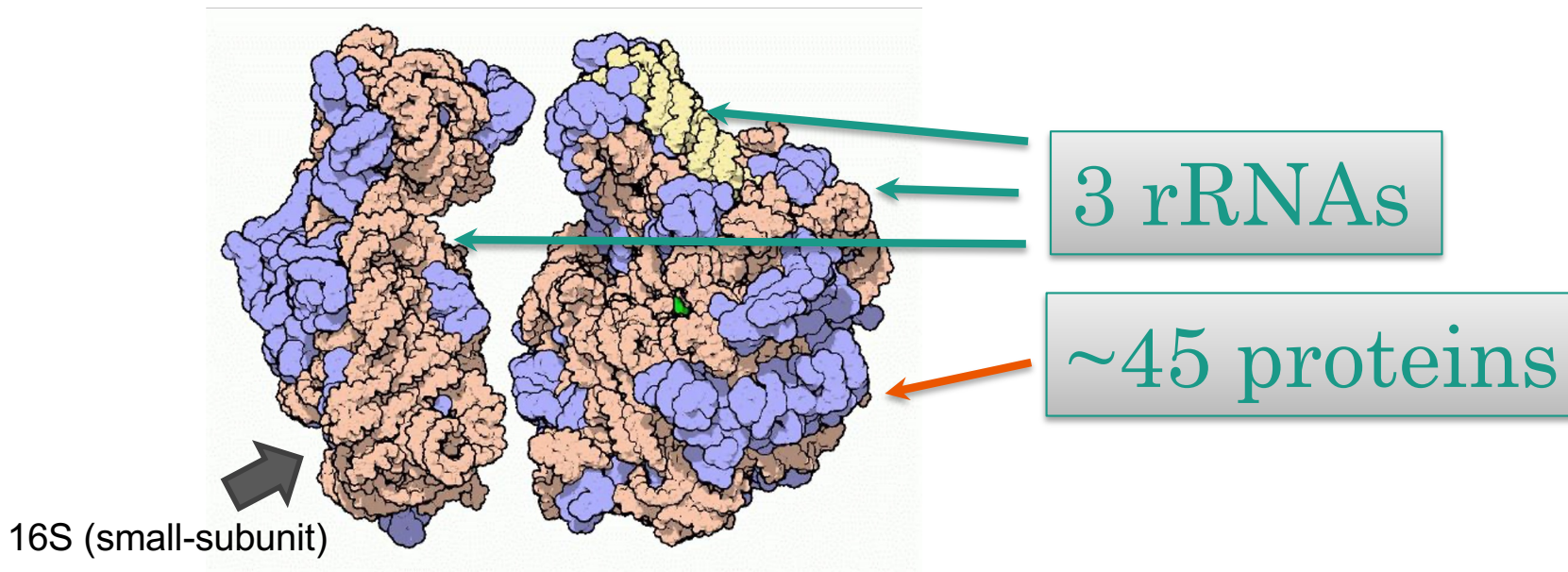
Taxonomy as a genomic game

- Great!
 - Lots of “objective” characters to compare (e.g., nucleotides, amino acids)
 - More phylogenetically informative than morphology
 - Can be assessed in an automated way
 - Less prone to convergence, maybe
- But:
 - We cannot profile the complete genomes of all organisms in any but the simplest of microbial communities
 - So we need to choose a smaller sequence that everyone has...
 - ...with enough characters to distinguish different groups
 - Variation and taxonomy may not have much to do with function (so won't tell us about important evolutionary events)
 - Molecular evolution is *weird*

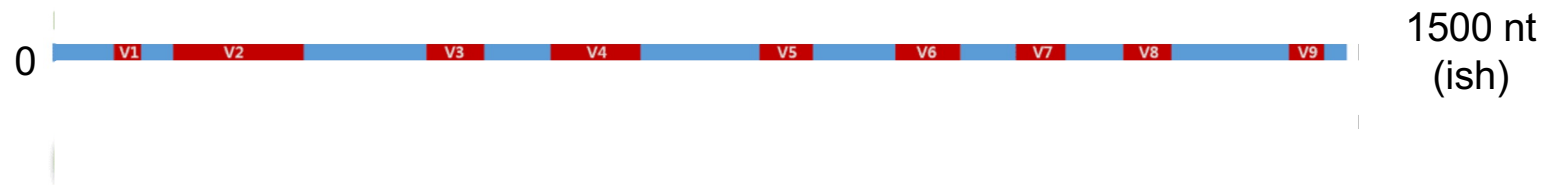
Genes that everyone has, more or less

- AMPHORA for bacterial core-genome phylogeny uses *dnaG*, *frr*, *infC*, *nusA*, *pgk*, *pyrG*, *rplA*, *rplB*, *rplC*, *rplD*, *rplE*, *rplF*, *rplK*, *rplL*, *rplM*, *rplN*, *rplP*, *rplS*, *rplT*, *rpmA*, *rpoB*, *rpsB*, *rpsC*, *rpsE*, *rpsI*, *rpsJ*, *rpsK*, *rpsM*, *rpsS*, *smpB*, *tsf*
- 23 of these encode ribosomal proteins!
- Which are great, but there are...issues with using protein-coding genes.
- So...

The prokaryotic ribosome (more or less)

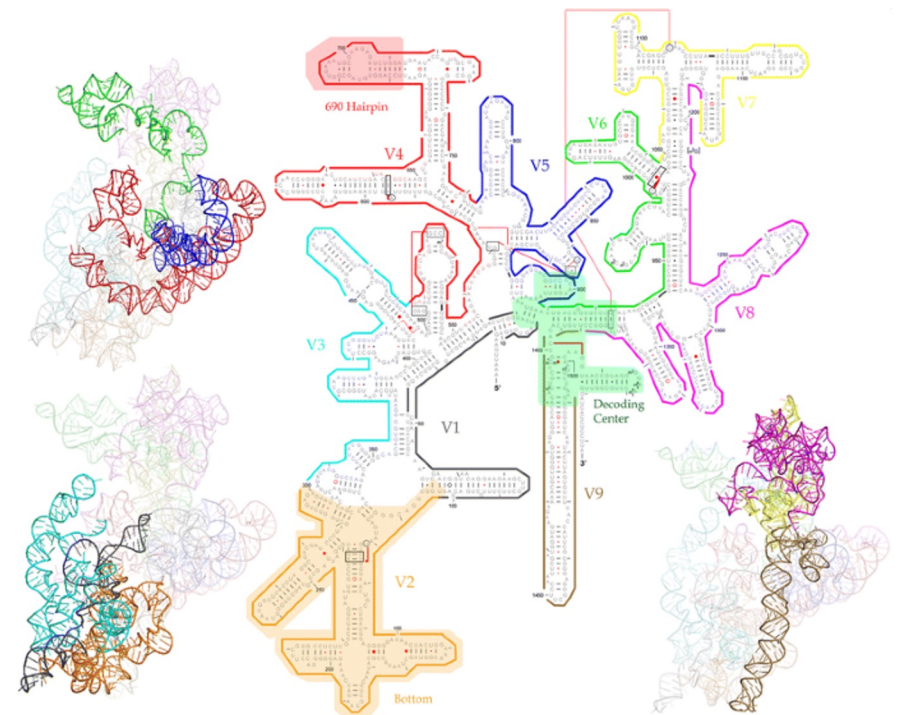


The 16S rRNA gene



Modified from <https://help.ezbiocloud.net/16s-rrna-and-16s-rrna-gene/>

Not just a straight line!





The Story of 16S

- Good!
 - Everyone has it
 - Conserved and variable regions, none of this genetic code nonsense
 - An absurd number of sequences can be found in public databases (9,469,124 in SILVA 138.1)
- Less good!
 - Found in multiple copies, recombination and gene transfer are things
 - Insufficient resolution at the finest (strain) levels
 - Differential amplification of groups depending on primer affinity (not 16S's fault)

16S is not the only option!

OPEN ACCESS Freely available online



The Chaperonin-60 Universal Target Is a Barcode for Bacteria That Enables *De Novo* Assembly of Metagenomic Sequence Data

Matthew G. Links^{1,2}, Tim J. Dumonceaux^{1,2}, Sean M. Hemmingsen^{3,4}, Janet E. Hill^{2*}

1 Agriculture and AgriFood Canada, Saskatoon, Saskatchewan, Canada, **2** Department of Veterinary Microbiology, University of Saskatchewan, Saskatoon, Saskatchewan, Canada, **3** National Research Council Canada, Saskatoon, Saskatchewan, Canada, **4** Department of Microbiology and Immunology, University of Saskatchewan, Saskatoon, Saskatchewan, Canada

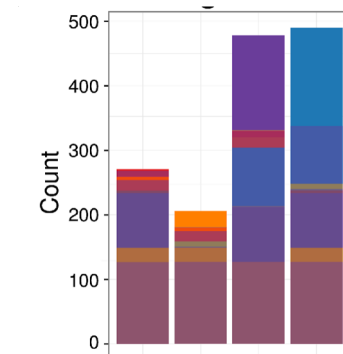
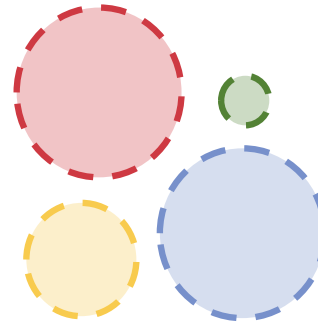
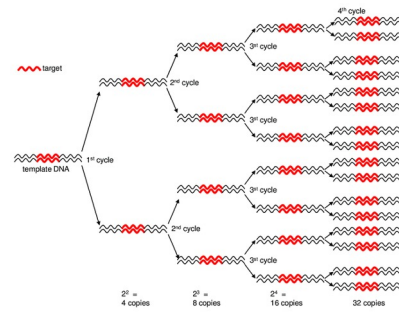
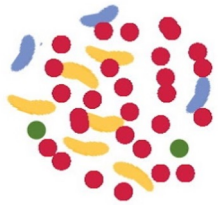
...

pyrosequencing data from a synthetic microbial community. Analysis supported the recognition of both 16S rRNA and *cpn60* as DNA barcodes for Bacteria. The *cpn60* universal target was found to have a much larger barcode gap than 16S rRNA suggesting *cpn60* as a preferred barcode for Bacteria. A large barcode gap for *cpn60* provided a robust target for species-level characterization of data. The assembly of consensus sequences for barcodes was shown to be a reliable method for the identification and tracking of novel microbes in metagenomic studies.

cpndb: <https://www.cpndb.ca/publications.php>

(2012)

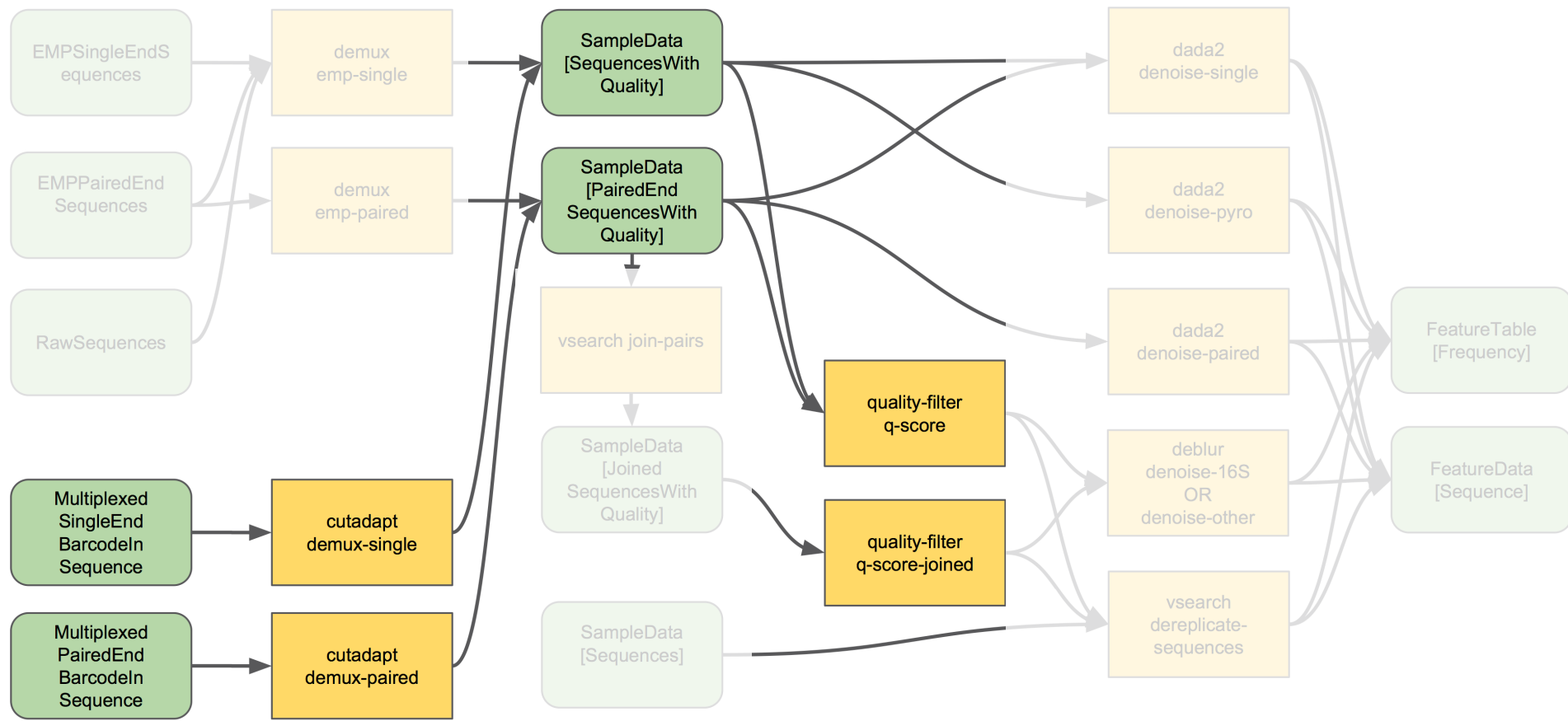
The basic microbial diversity pipeline



Sample → DNA → Amplify and sequence → Cluster → Diversity analysis

Plus statistics and possibly machine learning

Sequence retrieval and quality control





Anatomy of a FASTQ file

Header line

Sequence

+

PHRED Quality scores

@M03730:144:000000000-CG9WN:1:1101:16483:2296 1:N:0:377
ACGCGAAAACCTTACCAGGTCTTGACATCTAG...
+
CCCCDGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGGG...
→

'C': ASCII code 67,
PHRED score 34
(good!)



Importing into QIIME2

- Sequencing data will generally be in FASTQ format
- Paired-end reads will have two associated files, usually designated R1 and R2:

```
-rwxrwxrwx 1 rbeiko rbeiko 14678002 Mar 27 09:15 1051C-M-D00-V_S353_L001_R1_001.fastq.gz  
-rwxrwxrwx 1 rbeiko rbeiko 19940005 Mar 27 09:15 1051C-M-D00-V_S353_L001_R2_001.fastq.gz
```

- Files are zipped to save space; this is fine as you can still view contents (zmore) and QIIME2 works directly on the zipped files

Is it a real sequence?

I searched against the RefSeq database using the previous sequence as query (you generally won't do this). "Subject" is a 16S gene from *Lactobacillus plantarum*

Looks like it's legit...mismatches could be real, or sequencing errors

Lactobacillus crispatus strain 3019 16S ribosomal RNA gene, partial sequence

Sequence ID: [MT613437.1](#) Length: 1474 Number of Matches: 1

Range 1: 962 to 1262 [GenBank](#) [Graphics](#)

[▼ Next Match](#) [▲ Previous Match](#)

Score	Expect	Identities	Gaps	Strand
540 bits(292)	8e-150	298/301(99%)	0/301(0%)	Plus/Plus
Query 1	ACGCGAAAAACCTTACCAGGTCTTGACATCTAGTGCCATTTGTAGAGATACAAAGTTCCC	60		
Sbjct 962	ACGCGAAGAACCTTACCAGGTCTTGACATCTAGTGCCATTTGTAGAGATACAAAGTTCCC	1021		
Query 61	TTCGGGGACGCTAAGACAGGTGGTGCATGGCTGTCGTCAGCTCGTGTGAGATGTTGG	120		
Sbjct 1022	TTCGGGGACGCTAAGACAGGTGGTGCATGGCTGTCGTCAGCTCGTGTGAGATGTTGG	1081		
Query 121	GTTAAGTCCCGCAACGAGCGCAACCCCTGTTATTAGTTGCCAGCATTAAAGTTGGGCACTC	180		
Sbjct 1082	GTTAAGTCCCGCAACGAGCGCAACCCCTGTTATTAGTTGCCAGCATTAAAGTTGGGCACTC	1141		
Query 181	TAATGAGACTGCCGGTGACAAACGAGGAAAGTGGGGATGACGTC AAGTCATCATGCC	240		
Sbjct 1142	TAATGAGACTGCCGGTGACAAACGAGGAAAGTGGGGATGACGTC AAGTCATCATGCC	1201		
Query 241	CTTATGACCTGGGCTACACACGTGCTACAATGGGCGGACAAACGAGAAGCGAGCCTGCGA	300		
Sbjct 1202	CTTATGACCTGGGCTACACACGTGCTACAATGGGCGGACAAACGAGAAGCGAGCCTGCGA	1261		
Query 301	A 301			
Sbjct 1262	A 1262			



Importing into QIIME2

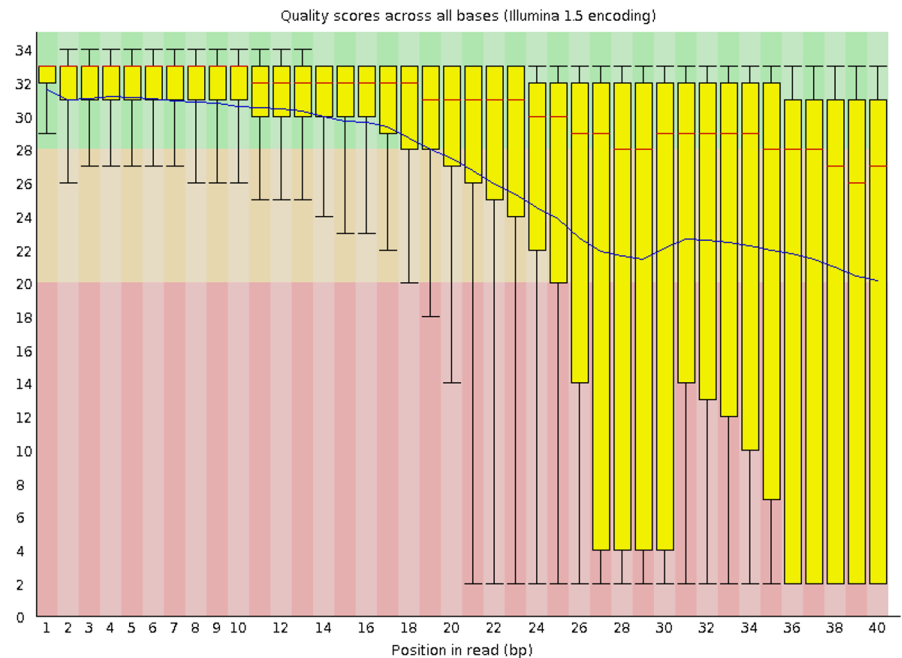
```
qiime tools import \  
--type 'SampleData[PairedEndSequencesWithQuality]' \  
--input-path import_to_qiime \  
--output-path CBW_reads
```

QIIME wants all the fastq.gz files to be in the same directory ('import_to_qiime' in this case)

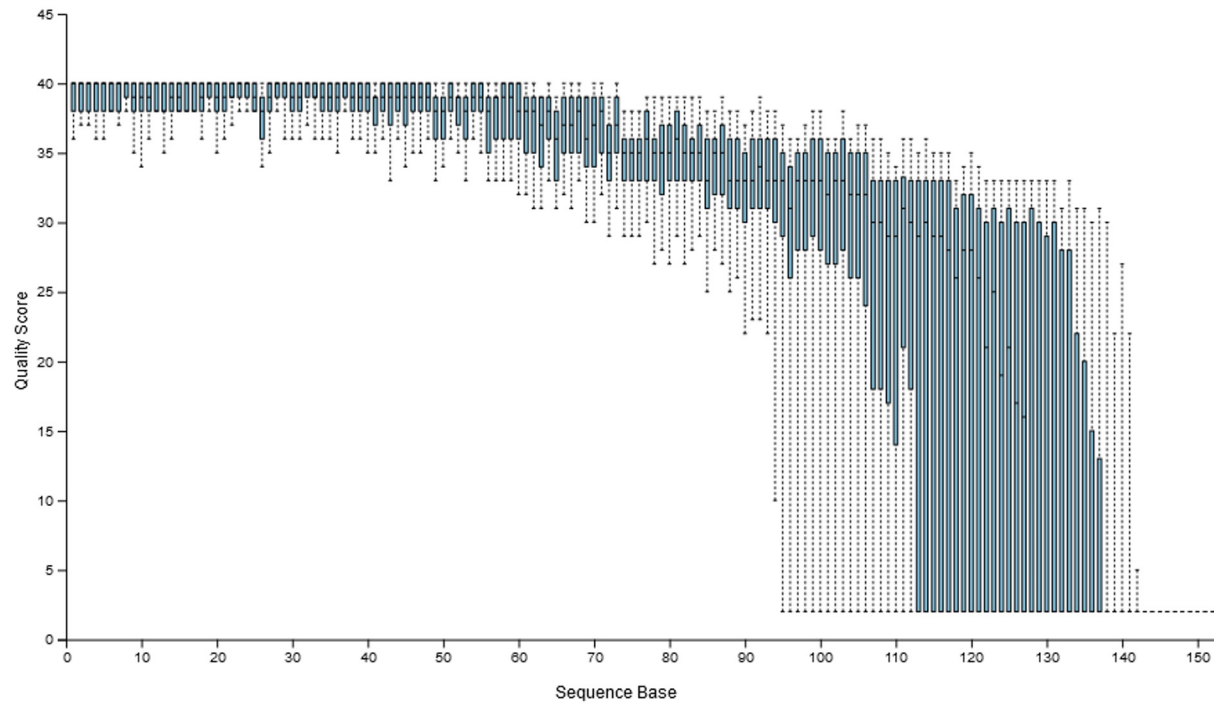
Visualizing sequence quality

Average over all reads in a dataset

Can show red flags of poor sequencing run



QIIME2 interactive quality plot (from “demux”)

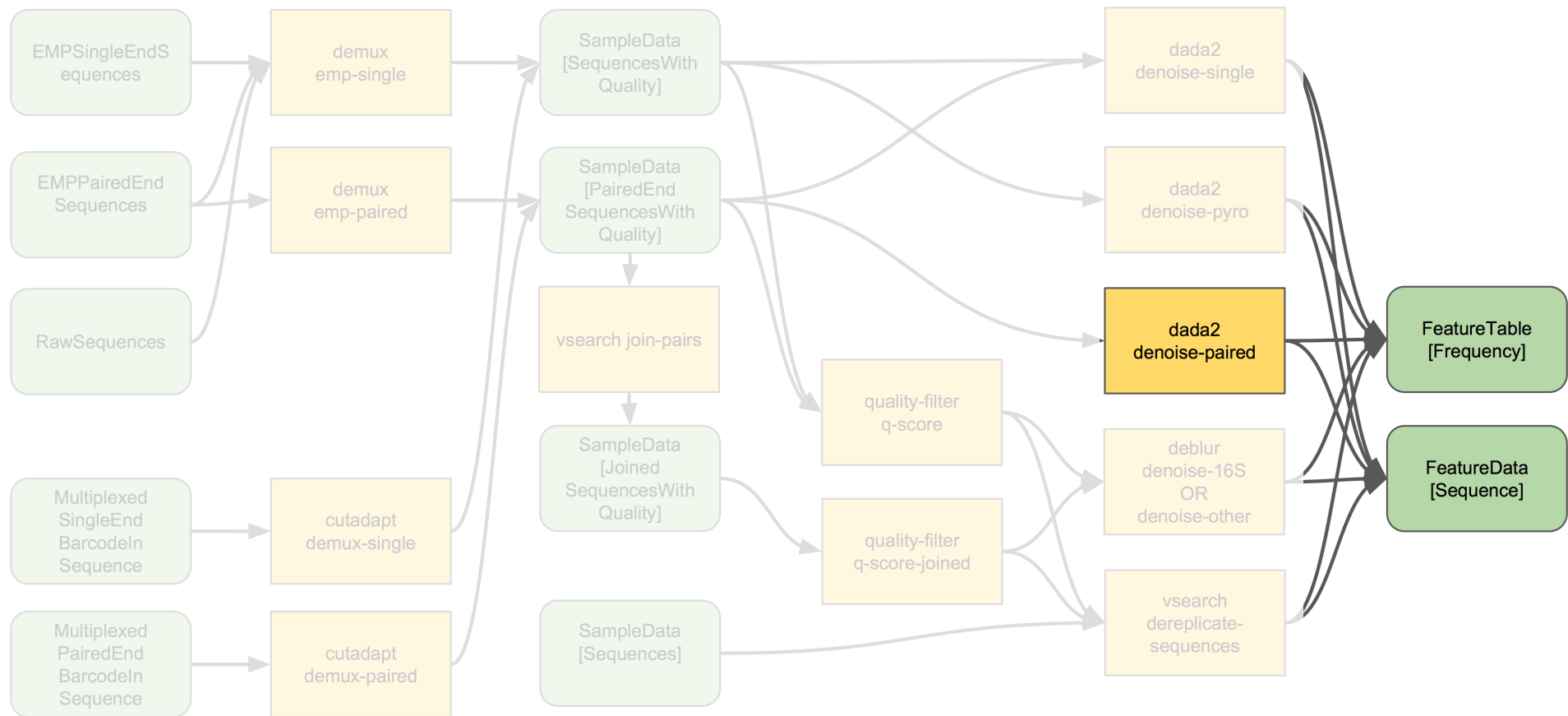


Files are imported. What now?

- If sequence files were multiplexed (i.e., multiple samples in a single file), 'cutadapt' can be used to assign sequences to samples according to provided barcodes.
- You can also remove reads below a certain length in 'cutadapt'
- Additional processing and use of quality scores will be carried out in the DADA2 step (next module)

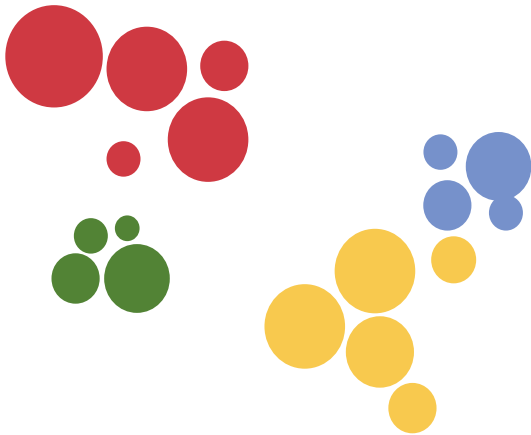
Sequence clustering

—

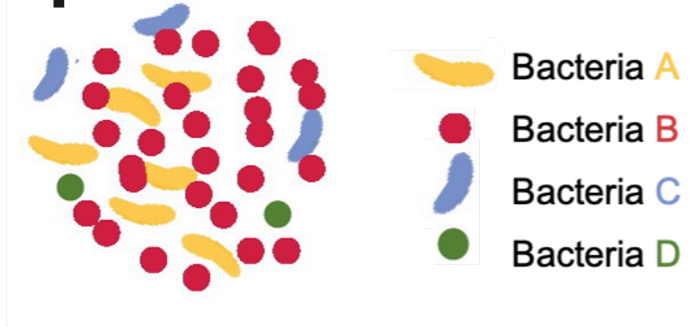


Why do we cluster sequences?

- Reduce dimensionality
- Remove sequencing artifacts
- Compensate for within-species variability



Operational Taxonomic Units (OTUs)



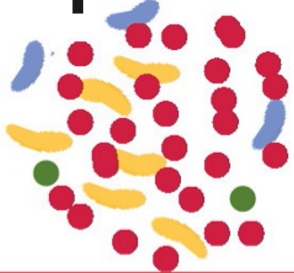
ATCGATCGATCGAT**T**GCTAG**A**TATCGATCGATCGATGCTAGCTATCGATCGATCGATGCTAGCTATCGA
ATCGATCGATCGA**G**GCTAG**C**TATCGATCGATCGATGCTAGCTATCGATCGATCGATGCTAGCTATCGA

$$\% \textit{ identity} = \frac{\# \textit{ matches}}{\textit{ length of sequence}} \cdot 100$$

$$\% \textit{ identity} = \frac{98}{100} \cdot 100$$

$$\% \textit{ identity} = 98\%$$

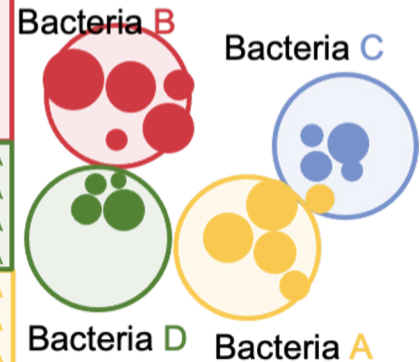
Operational Taxonomic Units (OTUs)



- Bacteria A
- Bacteria B
- Bacteria C
- Bacteria D

99% similarity threshold

ATCGATCGATCGATGCTAGATATCGATCGATCGATGCTAGCTATCGATCGATCGATGCTAGCTATCGATCGATCGATGCTAGCTGATCGATGCTAGCTA
ATCGATCGATCGAGGCTAGCTATCGATCGATCGATGCTAGCTATCGATCGATCGATGCTAGCTATCGATCGATCGATGCTAGCTGATCGATGCTAGCTA
ATCGATCGATCGATGCTAGCTATCGATCGATCGATGCTAGCTATCGATCGATCGATGCTAGCTATCGATCGATCGATGCTAGCTGATCGATGCTAGCTA
ATCGATCGATCGATGCTAGCTATCGATCGATCGATGCTAGCTATCGATCGATCGATGCTAGCTATCGATCGATCGATGCTAGCTGATCGATGCTAGCTA
ATCGATCGATCGATGCTAGCTATCGATCGATCGATGCTAGCTATCGATCGATCGATGCTAGCTATCGATCGATCGATGCTAGCTGATCGATGCTAGCTA
ATCGATCGATCGATGCTAGCTATCGATCGATCGATGCTAGCTATCGATCGATCGATGCTAGCTATCGATCGATCGATGCTAGCTGATCGATGCTAGCTA
ATCGATCGATCGATGCTAGCTATCGATCGATCGATGCTAGCTATCGATCGATCGATGCTAGCTATCGATCGATCGATGCTAGCTGATCGATGCTAGCTA
ATCGATCGATCGATGCTAGCTATCGATCGATCGATGCTAGCTATCGATCGATCGATGCTAGCTATCGATCGATCGATGCTAGCTGATCGATGCTAGCTA
ATCGATCGATCGATGCTAGCTATCGATCGATCGATGCTAGCTATCGATCGATCGATGCTAGCTATCGATCGATCGATGCTAGCTGATCGATGCTAGCTA
ATCGATCGATCGATGCTAGCTATCGATCGATCGATGCTAGCTATCGATCGATCGATGCTAGCTATCGATCGATCGATGCTAGCTGATCGATGCTAGCTA
ATCGATCGATCGATGCTAGCTATCGATCGATCGATGCTAGCTATCGATCGATCGATGCTAGCTATCGATCGATCGATGCTAGCTGATCGATGCTAGCTA
ATCGATCGATCGATGCTAGCTATCGATCGATCGATGCTAGCTATCGATCGATCGATGCTAGCTATCGATCGATCGATGCTAGCTGATCGATGCTAGCTA
ATCGATCGATCGATGCTAGCTATCGATCGATCGATGCTAGCTATCGATCGATCGATGCTAGCTATCGATCGATCGATGCTAGCTGATCGATGCTAGCTA
ATCGATCGATCGATGCTAGCTATCGATCGATCGATGCTAGCTATCGATCGATCGATGCTAGCTATCGATCGATCGATGCTAGCTGATCGATGCTAGCTA
ATCGATCGATCGATGCTAGCTATCGATCGATCGATGCTAGCTATCGATCGATCGATGCTAGCTATCGATCGATCGATGCTAGCTGATCGATGCTAGCTA
ATCGATCGATCGATGCTAGCTATCGATCGATCGATGCTAGCTATCGATCGATCGATGCTAGCTATCGATCGATCGATGCTAGCTGATCGATGCTAGCTA



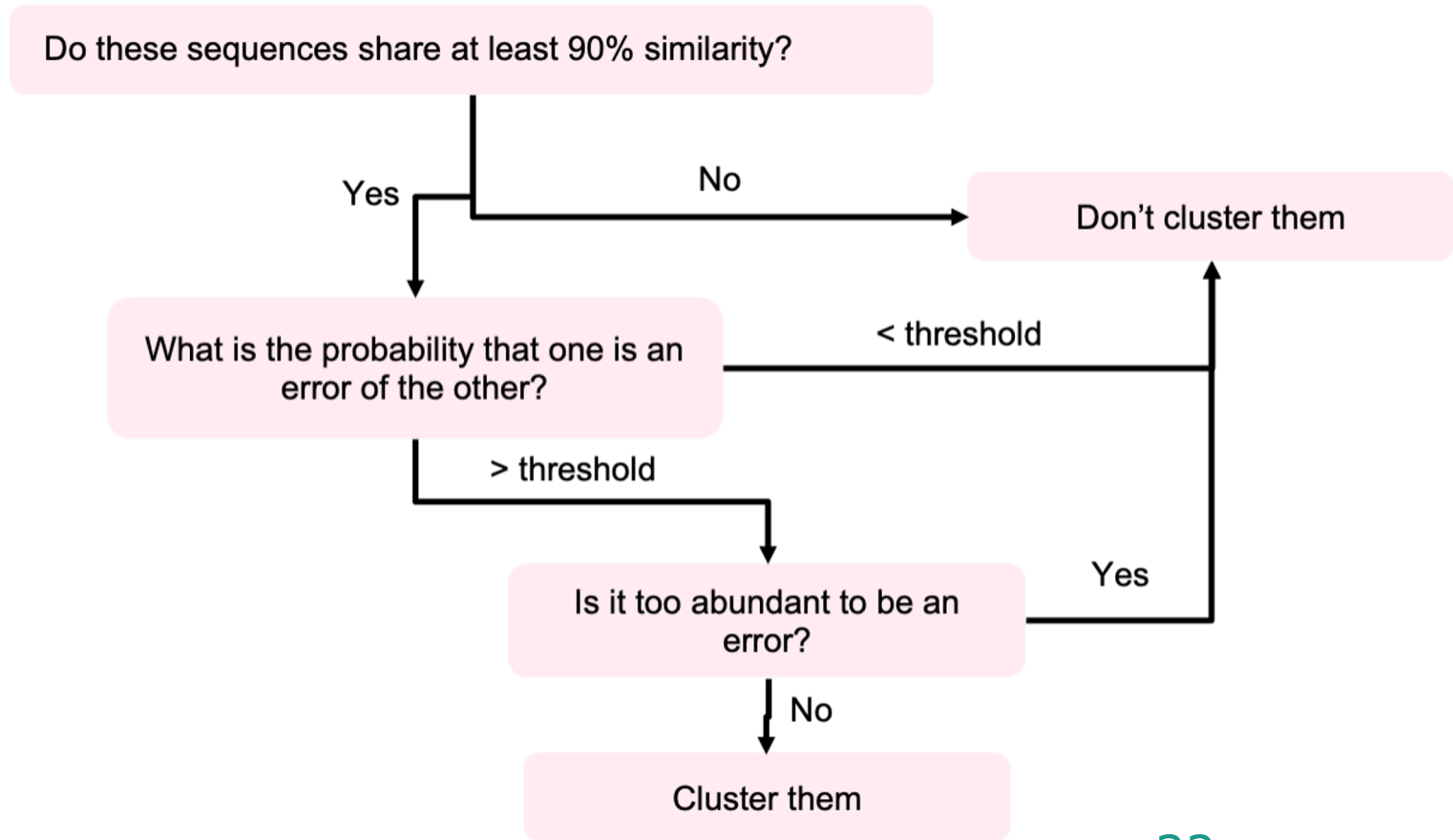
30

OTUs are a blunt instrument

- Grouping at arbitrary thresholds soaks up a lot of errors, but also loses a lot of information
- An alternative is calculating **Amplicon Sequence Variants*** that comprise “correct” sequences plus similar sequences with a high probability of sequencing errors

* or zero-radius OTUs, sub-OTUs

Inferring ASVs with DADA2

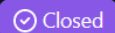


DADA2 1.10 RELEASE NOTES

NEW FEATURES

- PacBio CCS reads up to 3 kilobases are now supported. See also `PacBioErrFun`, the new and recommended error-estimation function for PacBio CCS data. The [preprint introducing DADA2's long-read functionality](#) has information on accuracy and sub-species resolution, and [the associated reproducible analyses show PacBio-specific workflows](#).

Can DADA2 pipeline process MinION whole 16S single read fastq files? #718

 Closed row2x opened this issue on Mar 27, 2019 · 2 comments



row2x commented on Mar 27, 2019

Is it possible to use DADA2 to process whole 16S single read fastq files generated by Nanoporetech MinION? The reads are on average 1500nt in length, and forward and reverse reads are interleaved in the same file. Thanks for your input.



benjjneb commented on Mar 27, 2019

Nope, not sensibly anyway. Nanopore error rates are way too high to meet DADA2's assumption that at least some complete error-free reads exist in the data.

Assignees

No one assigned

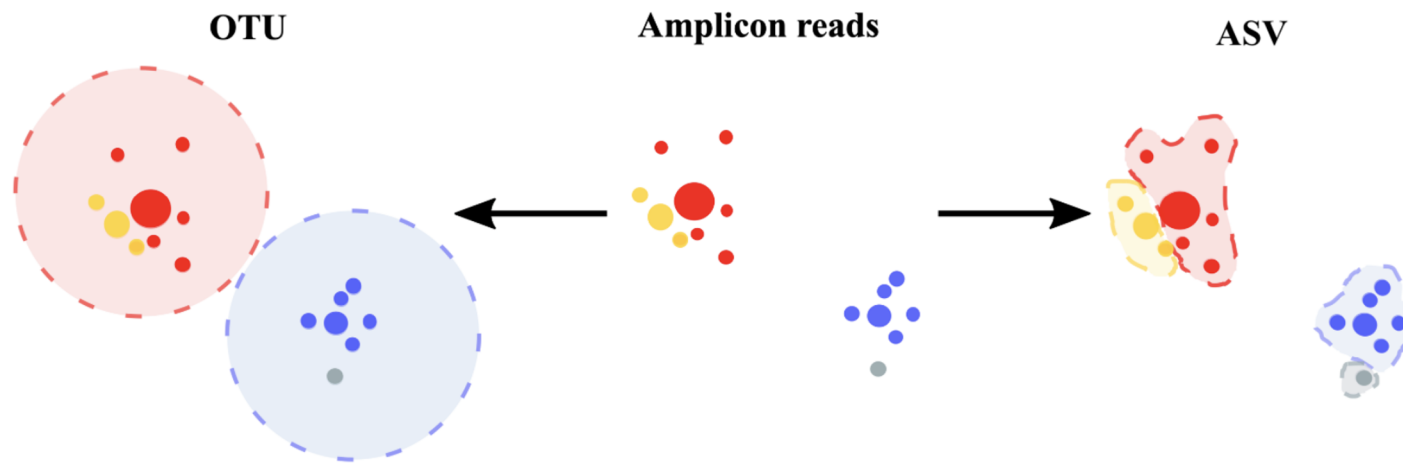
Labels

None yet

Projects

None yet

Milestone



- Alternatives
 - Deblur, UNOISE
 - All draw on notions of errors and relative abundance
 - Comparison: Nearing et al. (2018) *PeerJ*

How can I safely cluster my data?

- Try different methods,
- Understand them,
- Compare your results,
- Be able to defend your method



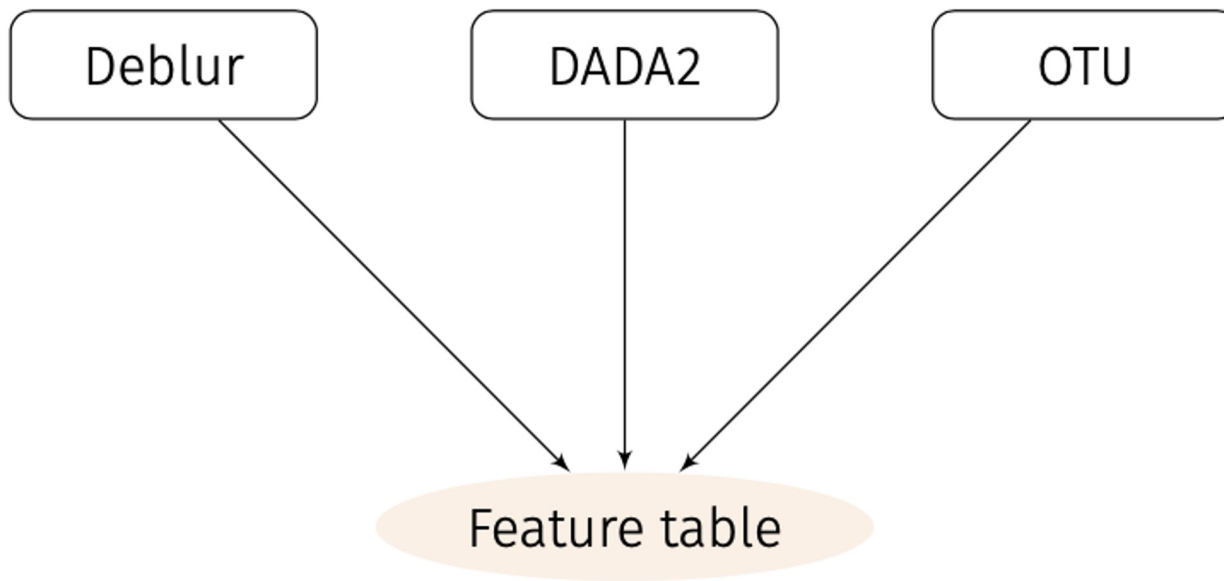


Table 1: Feature table

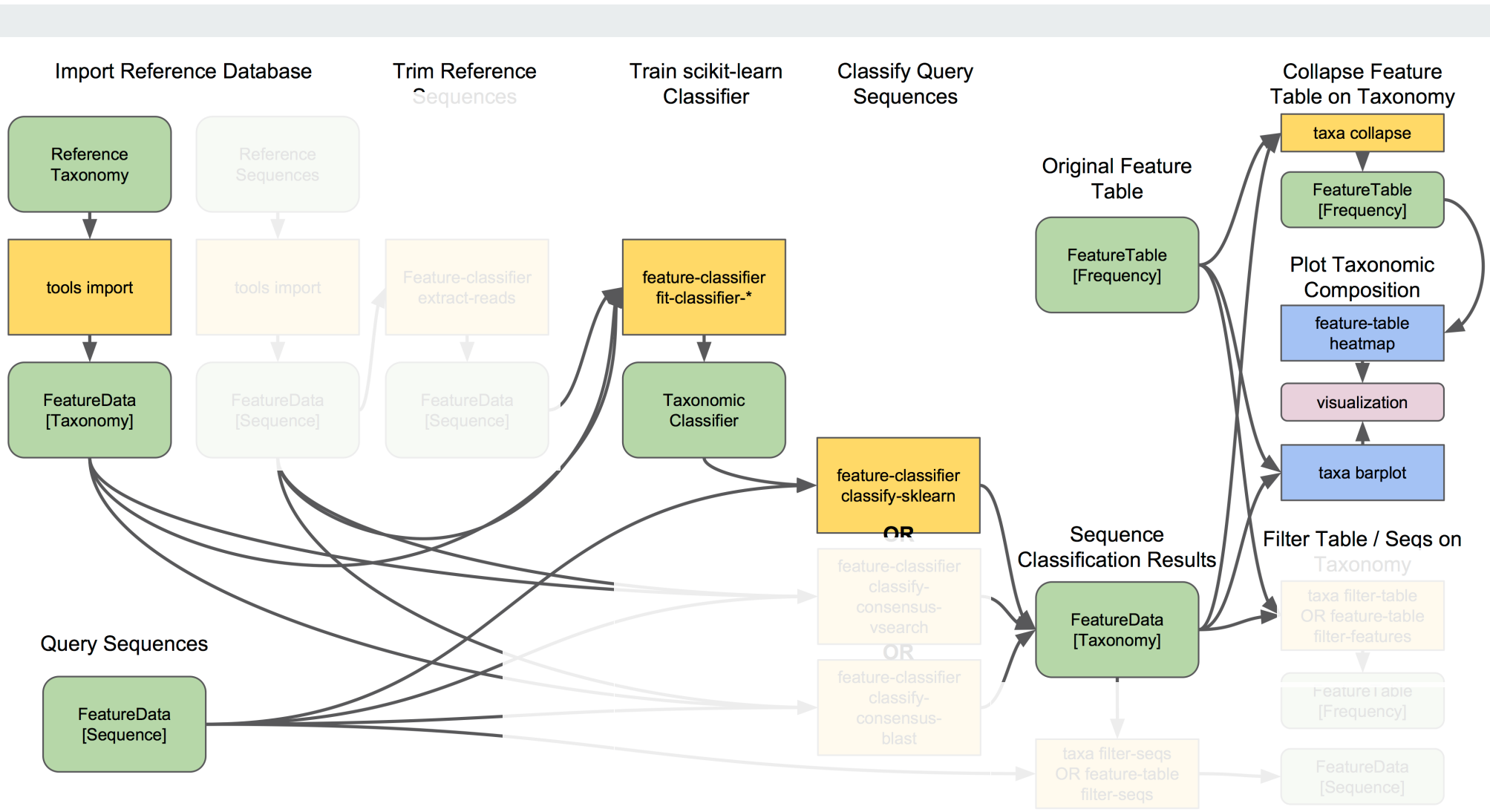
	Sample 1	Sample2	Sample 3	...
Microbe 1				
Microbe 2		Count		
...				

Taxonomic classification

From the feature table...

- Assign taxonomy to ASVs

feature_id	Taxon
00071b5a03322beb265e12561e62b7d2	d__Bacteria; p__Actinobacteriota; c__Acidimicrobiia; o__Actinomarinales; f__Actinomarinaceae; g__Candidatus_Actinoma
00071b5a03322beb265e12561e62b7d2	d__Bacteria; p__Actinobacteriota; c__Acidimicrobiia; o__Actinomarinales; f__Actinomarinaceae; g__Candidatus_Actinoma
0007c2bd79c09a2c5074f744573464e9	d__Bacteria; p__Proteobacteria; c__Gammaproteobacteria; o__Oceanospirillales; f__Pseudohongiellaceae; g__Pseudoho
0007c2bd79c09a2c5074f744573464e9	d__Bacteria; p__Proteobacteria; c__Gammaproteobacteria; o__Oceanospirillales; f__Pseudohongiellaceae; g__Pseudoho
0031efc0a916e54571c6ae4ddddbbd23	d__Bacteria; p__Proteobacteria; c__Gammaproteobacteria; o__Oceanospirillales; f__Pseudohongiellaceae; g__Pseudoho
0032e6d9726df6cf4176f545ca56623f	d__Bacteria; p__Proteobacteria; c__Alphaproteobacteria; o__Puniceispirillales; f__SAR116_clade; g__Candidatus_Punice
0032e6d9726df6cf4176f545ca56623f	d__Bacteria; p__Proteobacteria; c__Alphaproteobacteria; o__Puniceispirillales; f__SAR116_clade; g__Candidatus_Punice
003a1d54a176434c88b3e6488631bbe2	d__Bacteria; p__Proteobacteria; c__Gammaproteobacteria
003a1d54a176434c88b3e6488631bbe2	d__Bacteria; p__Proteobacteria; c__Gammaproteobacteria



Data resources

- Reference databases:
 - SILVA (updated 2021) – classifier in QIIME2
 - Greengenes (updated 2013) – classifier in QIIME2
 - RDP (updated 2016; taxonomy updated 2020)
 - EZBioCloud (updated 2021) – offers taxonomic profiling
 - NCBI (updated continuously)



Pat Schloss
@PatSchloss

Replying to @KevinDKohl

greengenes is basically dead. RDP uses Bergey's outline. SILVA is the way to go for alignment. Could flip a coin on RDP vs. SILVA for classification

5:44 PM · Apr 5, 2019 · Twitter Web Client



RESCRIPt: a tool for customized database retrieval and management:
<https://www.biorxiv.org/content/10.1101/2020.10.05.326504v1.full>

Taxonomic inference

- BLAST assigns taxonomic label of closest match
- Supervised learning algorithm : Naïve Bayes classifier

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$

↑ THE PROBABILITY OF "A" BEING TRUE GIVEN THAT "B" IS TRUE

↓ THE PROBABILITY OF "B" BEING TRUE GIVEN THAT "A" IS TRUE

↑ THE PROBABILITY OF "A" BEING TRUE

↑ THE PROBABILITY OF "B" BEING TRUE

Train your own, or use a pre-trained classifier

- The classifier's accuracy increases if the classifier is trained on the region of the rRNA that was amplified (eg. V3-V4 or V5-V6 of the 16S rRNA)
- But! If you use common marker genes (16S, 18S, ITS) There are a lot of pre-trained classifiers online trained on the full length gene, or on specific primers

Naive Bayes classifiers trained on:

- [Silva 138 99% OTUs full-length sequences](#) (MD5: b8609f23e9b17bd4a1321a8971303310)
- [Silva 138 99% OTUs from 515F/806R region of sequences](#) (MD5: e05afad0fe87542704be96ff483824d4)
- [Greengenes 13_8 99% OTUs full-length sequences](#) (MD5: 6bbc9b3f2f9b51d663063a7979dd95f1)
- [Greengenes 13_8 99% OTUs from 515F/806R region of sequences](#) (MD5: 9e82e8969303b3a86ac941ceafeeac86)

Filtering on “confidence”

- Each taxonomic classification from the feature table comes with an associated “confidence” score between 0 and 1.0
- In the default taxonomic assignment scheme, this is a posterior probability – don’t confuse “confidence” with any formal statistical notion such as “confidence interval”
- Anything with confidence below a threshold (default = 0.7) is not classified
 - Not classified at species level? Try classifying at genus instead, etc.

Bokulich et al. *Microbiome* (2018) 6:90
<https://doi.org/10.1186/s40168-018-0470-z>

Microbiome

RESEARCH

Open Access

Optimizing taxonomic classification of marker-gene amplicon sequences with QIIME 2’s q2-feature-classifier plugin



Nicholas A. Bokulich^{1†}, Benjamin D. Kaehler^{2†}, Jai Ram Rideout¹, Matthew Dillon¹, Evan Bolyen¹, Rob Knight³, Gavin A. Huttley^{2*} and J. Gregory Caporaso^{1,4*}

So?

- Be aware that there are many tools, many classifiers, many databases!
- Different feature tables can potentially give you very different outcomes in taxonomy, diversity, statistics, and ML



End of
46 Part II