



H3ABioNet

Pan African Bioinformatics Network for H3Africa

16SrRNA Intermediate Bioinformatics Online Course: Int_BT

Module 4

16S rRNA Sequencing Bioinformatics Pipelines: the theory



H3ABioNet

Pan African Bioinformatics Network for H3Africa



Learning Objectives

1. To understand various 16S rRNA Sequencing platforms
2. Understand how to assess the quality of raw sequences
3. To know various Bioinformatics Pipelines used to process 16S rRNA Microbiome data

Learning Outcomes

1. Use Appropriate Microbiome Sequencing Tools
2. Acquire Knowledge to Assess Quality of Raw Sequence Reads
3. Knowledge on the Available 16S rRNA Bioinformatics Processing Pipelines

Outline

1. Background on 16S rRNA gene
2. Sequencing Tools
3. Quality Assessment of Raw Sequence reads
4. 16S rRNA Bioinformatic processing Pipeline
 - i. UCT-CBIO microbiome processing pipeline
 - ii. QIIME2 microbiome processing pipeline
 - iii. MOTHUR microbiome processing pipeline
 - iv. DADA2 microbiome processing pipeline

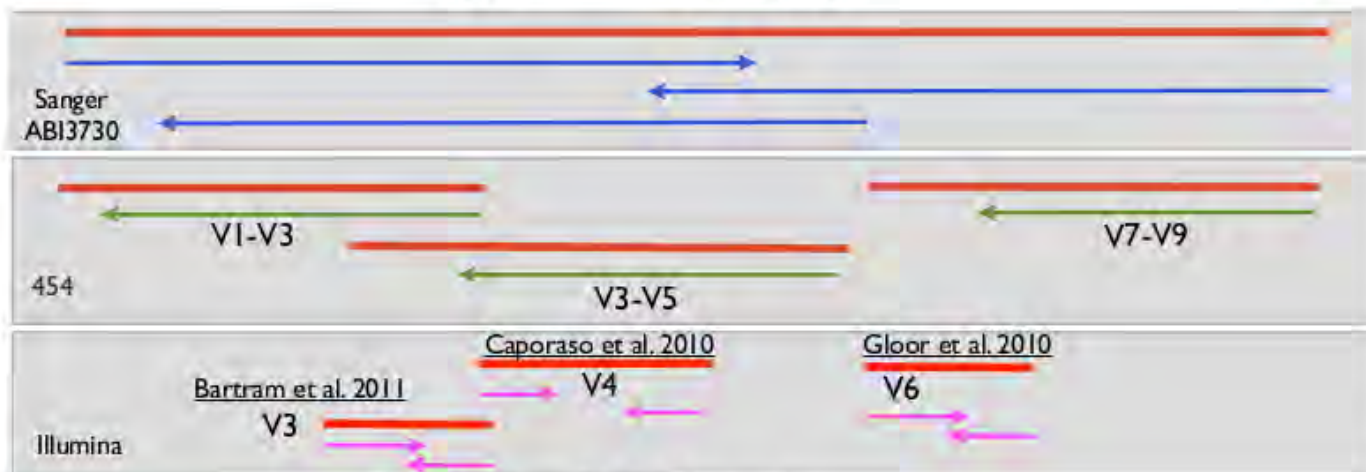
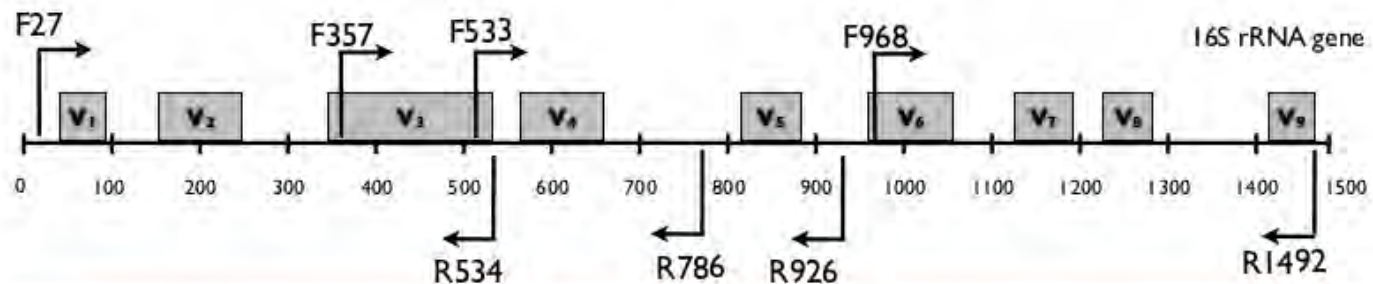
Module 4 : Part II

Background on 16S rRNA gene & Sequencing Platforms

Background on 16S rRNA

- The 16S rRNA is a gene which codes for the RNA component of the small subunit of ribosomes, and it is among the **most conserved** gene across all kingdoms of life.
- They contain **regions** that are less evolutionarily constrained and whose sequences are indicative of their phylogeny.
- **Amplification** of these genomic regions by PCR and subsequent **sequencing** of a sufficiently large number of individual **amplicons** enables the analysis of the **diversity** of organism in the sample and a rough estimate of their **relative abundance**.

16S rRNA gene region



	Read Length	Depth of Sequencing
Amplicon		
Sanger 3730 xl read	800-1000 bp	+
454 FLX/Titanium read	250-400 bp	+++
Illumina GA10k read	75-150 bp	+++++



H3ABioNet

Pan African Bioinformatics Network for H3Africa

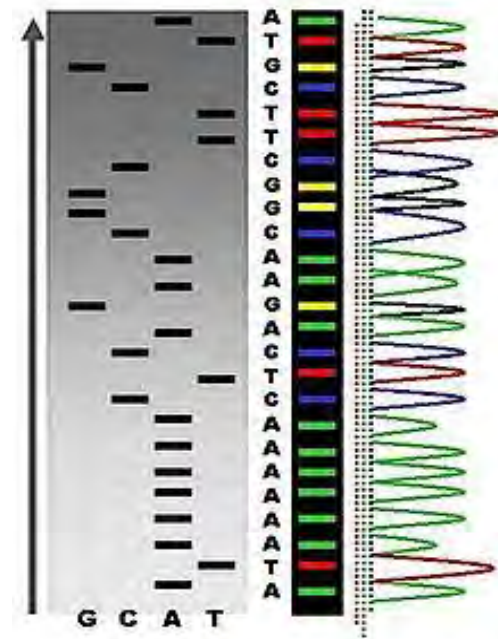
16SrRNA Intermediate Bioinformatics Online Course:

Int_BT_2019 Samson KM

Evolution of Sequencing platforms

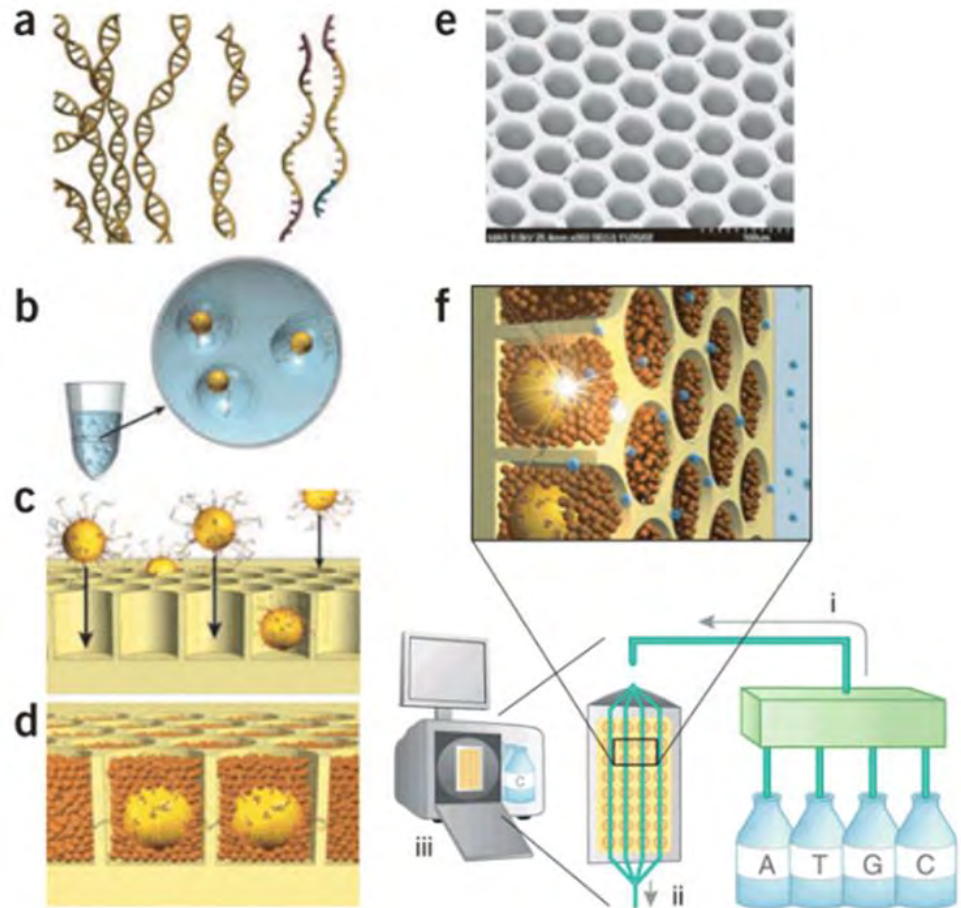
1. Sanger Sequencing

- Good quality
- Long reads, up to 1000 bp
- By comparison very little data
 - Up to 384 sequences read in parallel
- Expensive per base pair
- Developed by Frederick Sanger, 1977



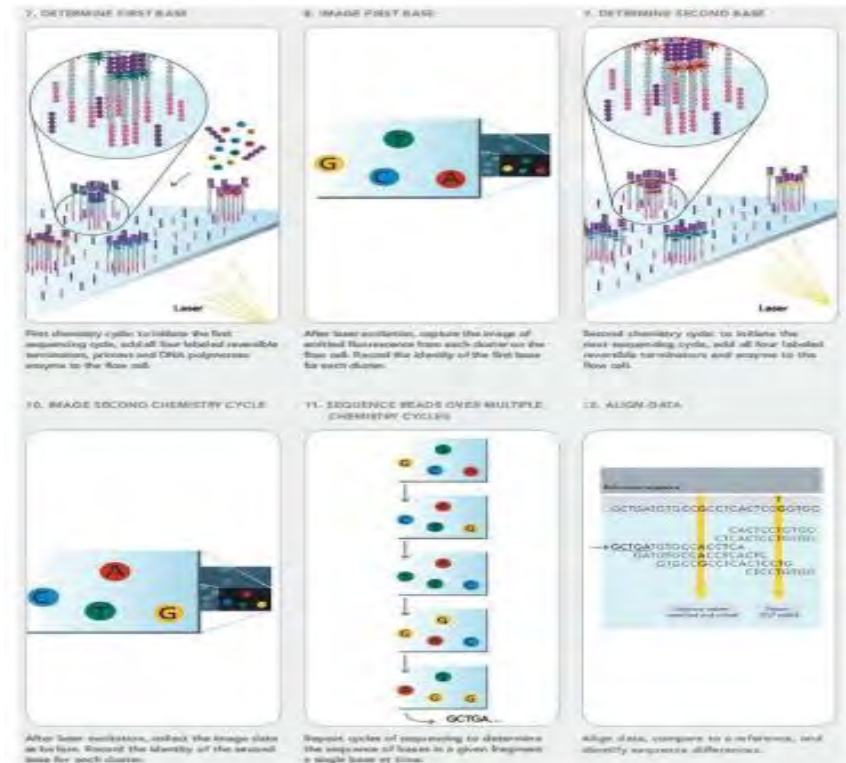
2. 454 and Ion Torrent

- Based on pyrosequencing
 - About 1 million sequences in parallel
 - 450bp, 700bp or longer
 - Expensive chemicals
 - More errors compared with Sanger
 - Homopolymers
 - Ion Torrent is similar
 - Shorter reads (100bp-400bp)
 - Much less expensive



3. Illumina

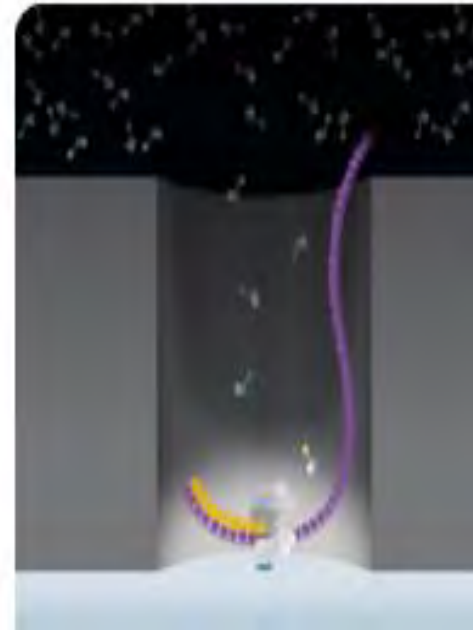
- Quite short reads: 2x80-150bp (HiSeq), 2x300bp (MiSeq)
- 3 billion reads per flow cell (8 lanes; HiSeq), 20 million reads (MiSeq)
- Much cheaper than 454
- More error compared both with Sanger and 454
 - Error distribution different to 454, single base pairs more frequent, homopolymer errors rare



From: <https://www.illumina.com/>

4. PacBio SMRT

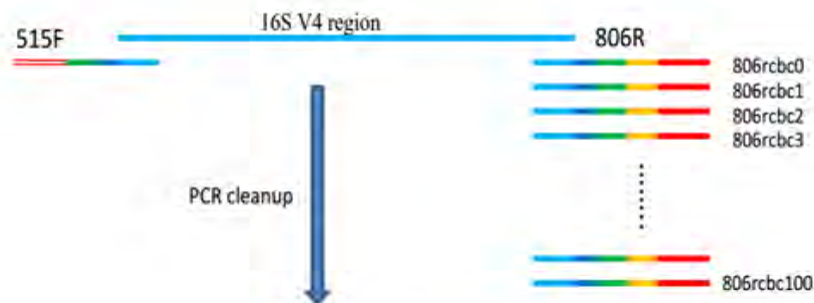
- Single Molecule Real-Time Sequencing nanotechnology:
 - Observe a single nucleotide being added to polymer (fluorescence)
- Long reads: 3 kb average
- 90 Mb per run
- Often combined with Illumina



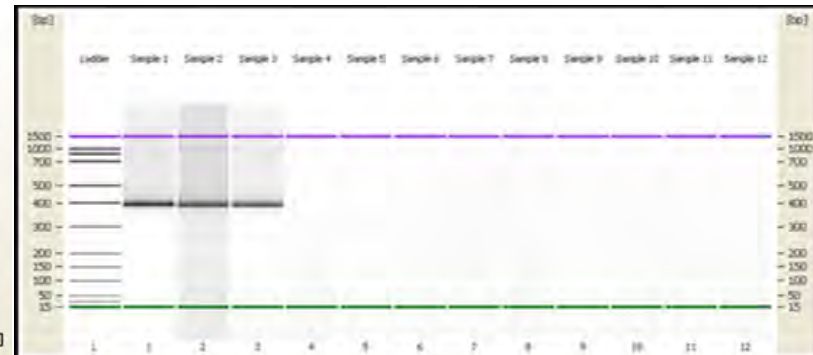
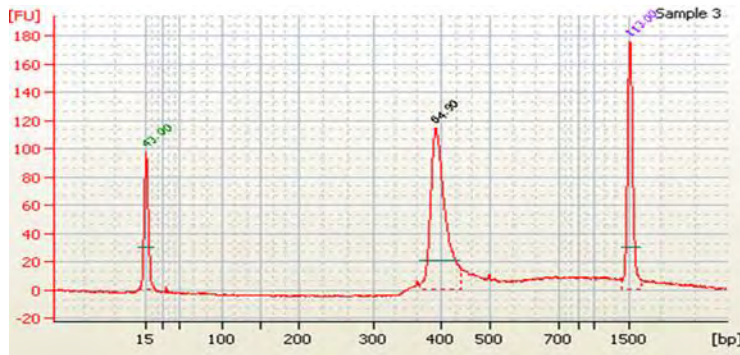
From: <http://www.pacificbiosciences.com>

SUMMARY of PCR and Sequencing Steps

Stage1 PCR

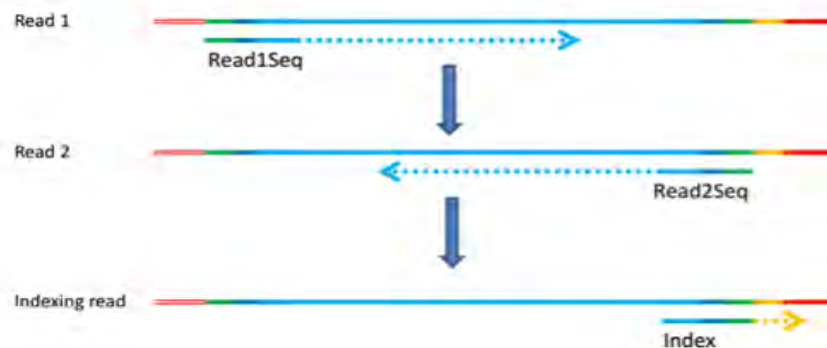


Stage 2 QC analysis



Stage 3 Sequencing

- P5 region, complementary to P5 region on Flowcell
- P7 region, complementary to P7 region on Flowcell
- Indices



Output file formats

- **SFF** (standard flowgram format) - 454
- **Fastq** - Illumina
- **BAM** (binary of sequence alignment map) – Ion Torrent
- **Metadata** in tabular format that can be used for downstream analysis

Module 4 : Part III


Quality Assessment of Raw Sequences

“Garbage in garbage out”

It takes a good lab practice to produce reliable data for the downstream processing

If we mess-up in Wetlab it can not be corrected in Dry lab

Assessment of the Quality of Raw Reads

 BM-AO-KITdnam-extractioncontrolspike-1-kitcomparison-P3-B07_S211_L001_R1_001.fastq
 BM-AO-KITdnam-extractioncontrolspike-1-kitcomparison-P3-B07_S211_L001_R2_001.fastq
 BM-AO-KITdnam-sequencing-control-1-wrrepeat-P3-C06_S222_L001_R1_001.fastq
 BM-AO-KITdnam-sequencing-control-1-wrrepeat-P3-C06_S222_L001_R2_001.fastq

```
@SEQ_ID
GATTGTTGGGGTTCAAAGCAGTATCGATCAAATAGTAAATCCATTTGTTCAACTCACAGTTT
+
!''*((( (**+)) %%%++) (%%%) .1***-+*'') **55CCF>>>>>CCCCCCC65
```

- Line 1 begins with a '@' character and is followed by a sequence identifier and an optional description (like a FASTA title line).
- Line 2 is the raw sequence letters.
- Line 3 begins with a '+' character and is optionally followed by the same sequence identifier (and any description) again.
- Line 4 encodes the quality values for the sequence in Line 2, and must contain the same number of symbols as letters in the sequence



Phred score

Phred score (Q) describes the quality of the sequences and is presented by an integer value. The larger the value the more confidence there can be in the output.

The probability a base is called incorrectly is given by

$$(P) = 10^{(-Q/10)}$$

Thus, $Q = -10 \log_{10} (P)$

- For example, if the probability of an error (P) = 0.01, then Q will be = 20. What if P = 0.0004? What will be the Q score?

ASCII_BASE=33 Illumina, Ion Torrent, PacBio and Sanger

Q	P_error	ASCII	Q	P_error	ASCII	Q	P_error	ASCII
0	1.00000	33 !	11	0.07943	44 ,	22	0.00631	55 7
1	0.79433	34 "	12	0.06310	45 -	23	0.00501	56 8
2	0.63096	35 #	13	0.05012	46 .	24	0.00398	57 9
3	0.50119	36 \$	14	0.03981	47 /	25	0.00316	58 :
4	0.39811	37 %	15	0.03162	48 0	26	0.00251	59 ;
5	0.31623	38 &	16	0.02512	49 1	27	0.00200	60 <
6	0.25119	39 '	17	0.01995	50 2	28	0.00158	61 =
7	0.19953	40 (18	0.01585	51 3	29	0.00126	62 >
8	0.15849	41)	19	0.01259	52 4	30	0.00100	63 ?
9	0.12589	42 *	20	0.01000	53 5	31	0.00079	64 @
10	0.10000	43 +	21	0.00794	54 6	32	0.00063	65 A



It is RECOMMENDED to check the quality of raw reads before further processing of the raw Sequences.

FASTQC tool by Andrew, 2010 ; is a very powerful open source quality control tool for assessing high throughput sequence data

Details found here:

<http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>

Example of good Illumina data:

https://www.bioinformatics.babraham.ac.uk/projects/fastqc/good_sequence_short_fastqc.html

Example of bad Illumina data

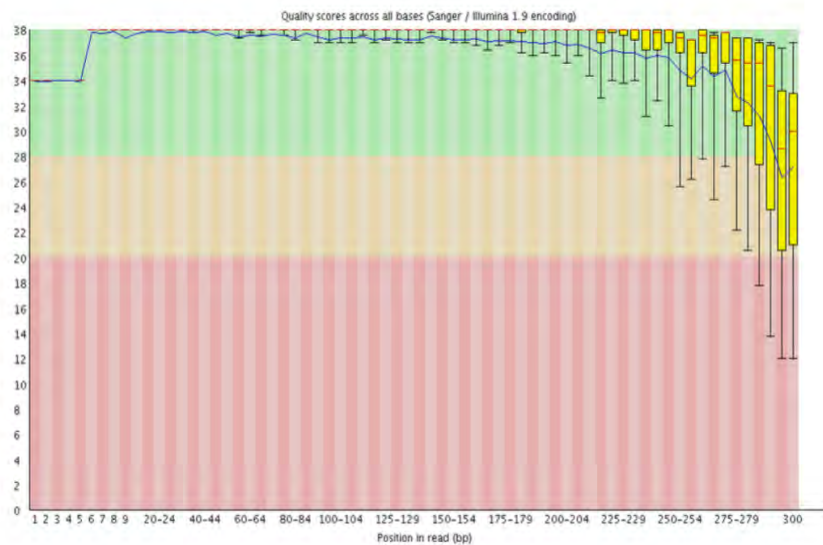
https://www.bioinformatics.babraham.ac.uk/projects/fastqc/bad_sequence_fastqc.html

Sample of FASTQC output (1)

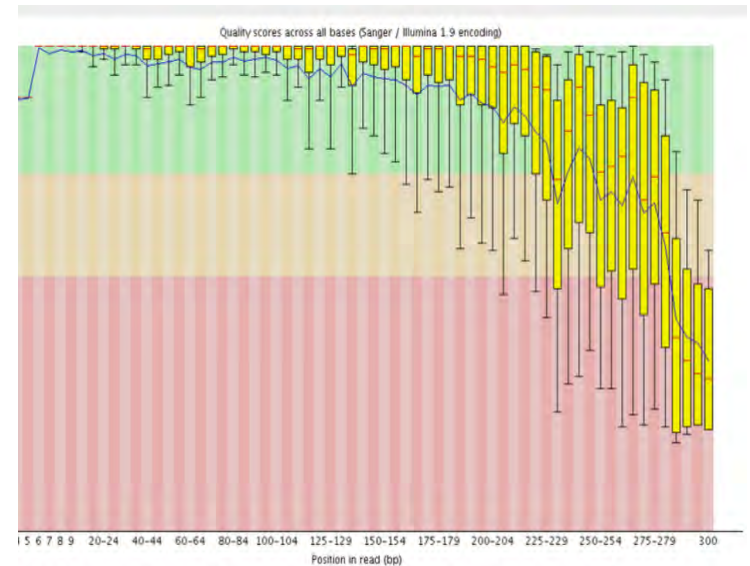
FastQC		
File Help		
IS-1703-case-230-88-516934-episode1event-R13-P1-A12_S12_L001_R2_001.fastq.gz		
IS-1703-case-230-88-516934-episode1event-R13-P1-A12_S12_L001_R1_001.fastq.gz		
Basic Statistics	Basic sequence stats	
	Measure	Value
Per base sequence quality	Filename	IS-1703-case-230-88-516934-episode1event-R13-P
Per tile sequence quality	File type	Conventional base calls
	Encoding	Sanger / Illumina 1.9
Per sequence quality scores	Total Sequences	19966
Per base sequence content	Sequences flagged as poor quality	0
	Sequence length	251
Per sequence GC content	%GC	51
Per base N content		
Sequence Length Distribution		
Sequence Duplication Levels		
Overrepresented sequences		
Adapter Content		

Sample FASQC output (2)

R1 (Forward read)

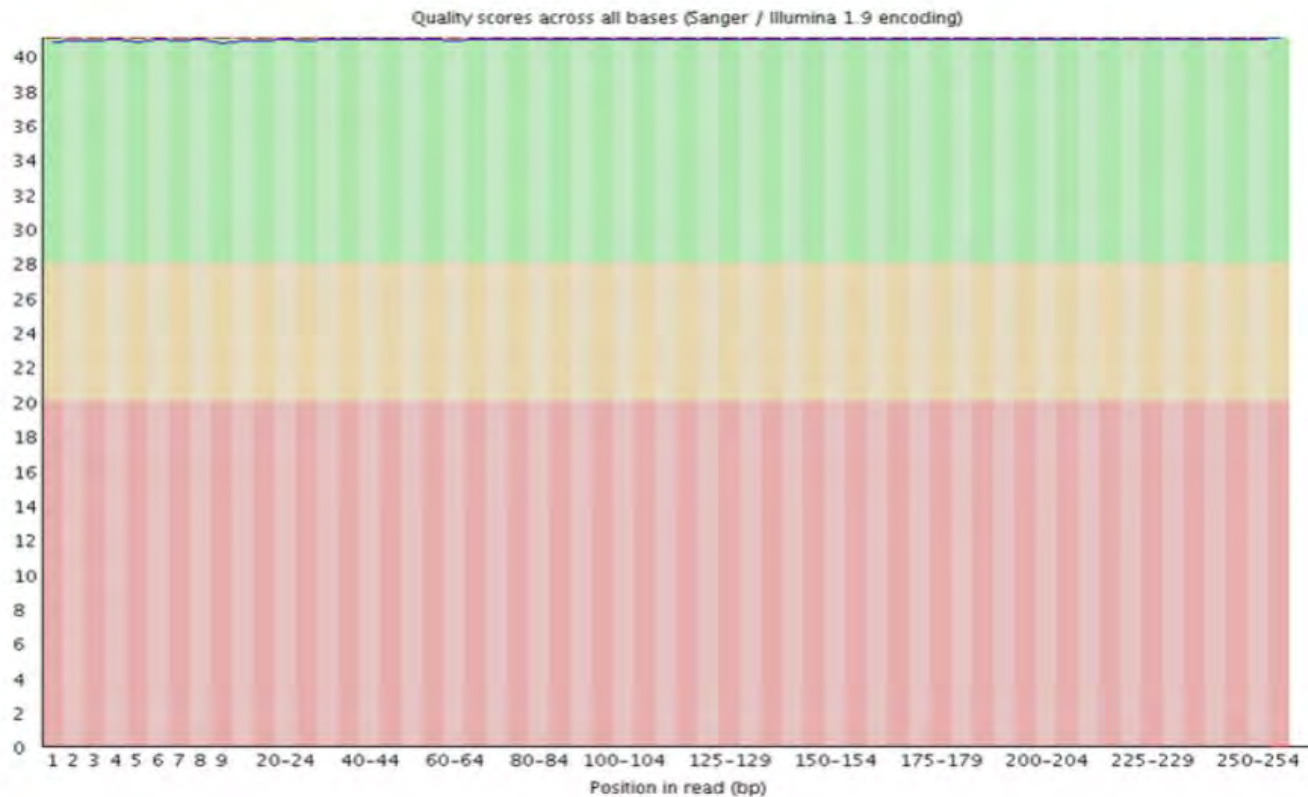


R2 (Reverse read)



What about this sample?

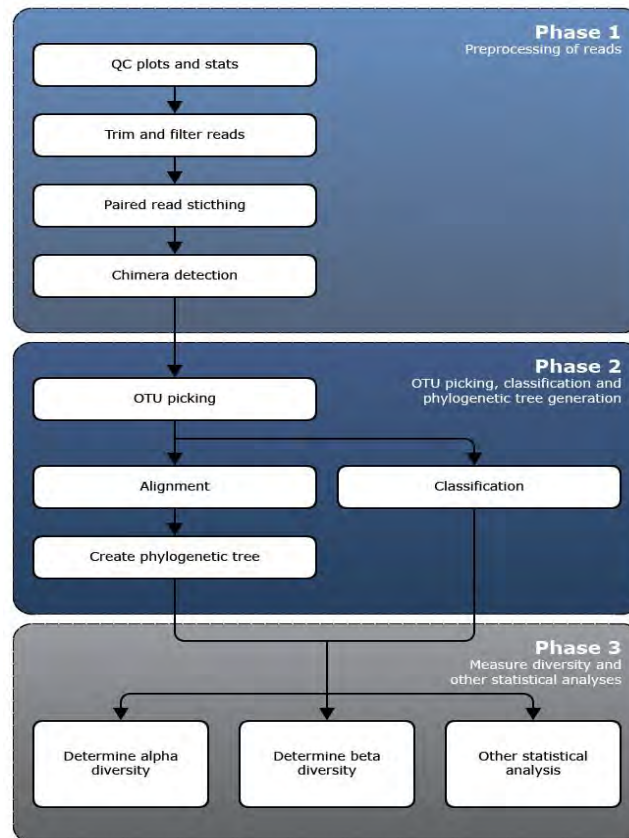
What about this?



Module 4: Part IV (a)

16S rRNA Bioinformatics Processing Pipelines

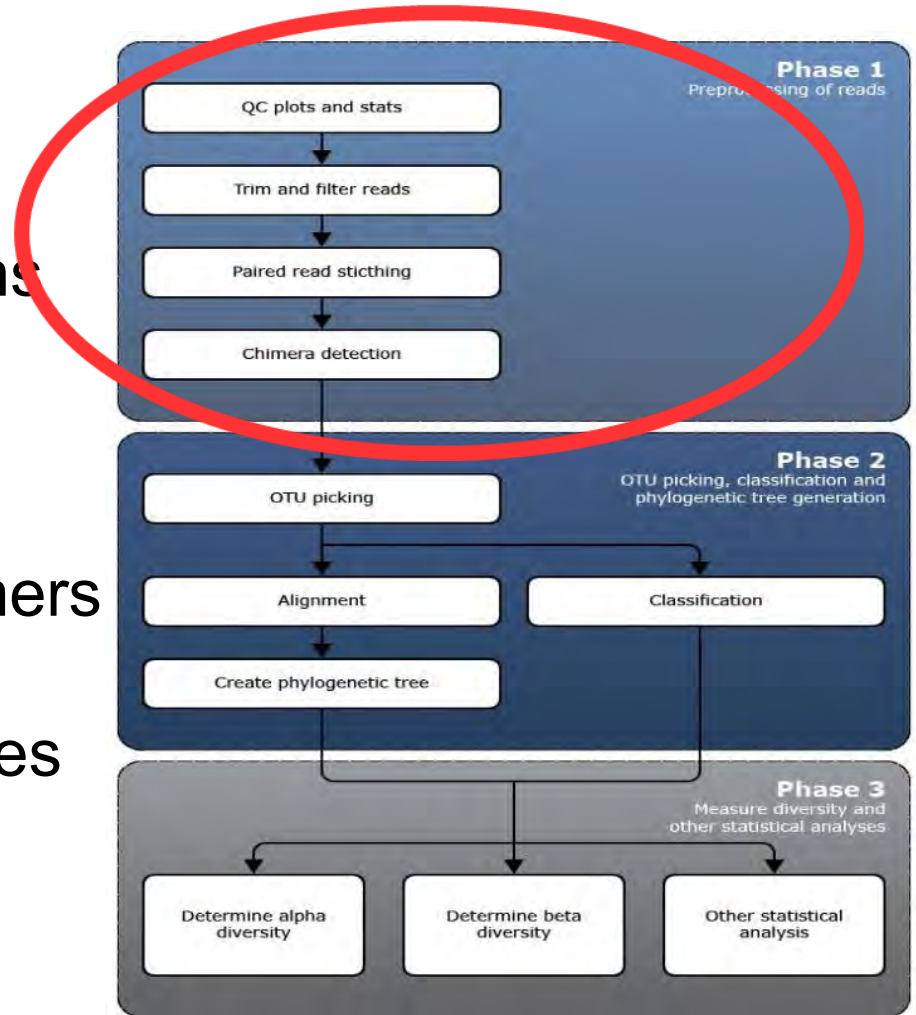
16S rRNA diversity analysis workflow



From: <http://h3abionet.org/tools-and-resources/sops/16s-rrna-diversity-analysis>

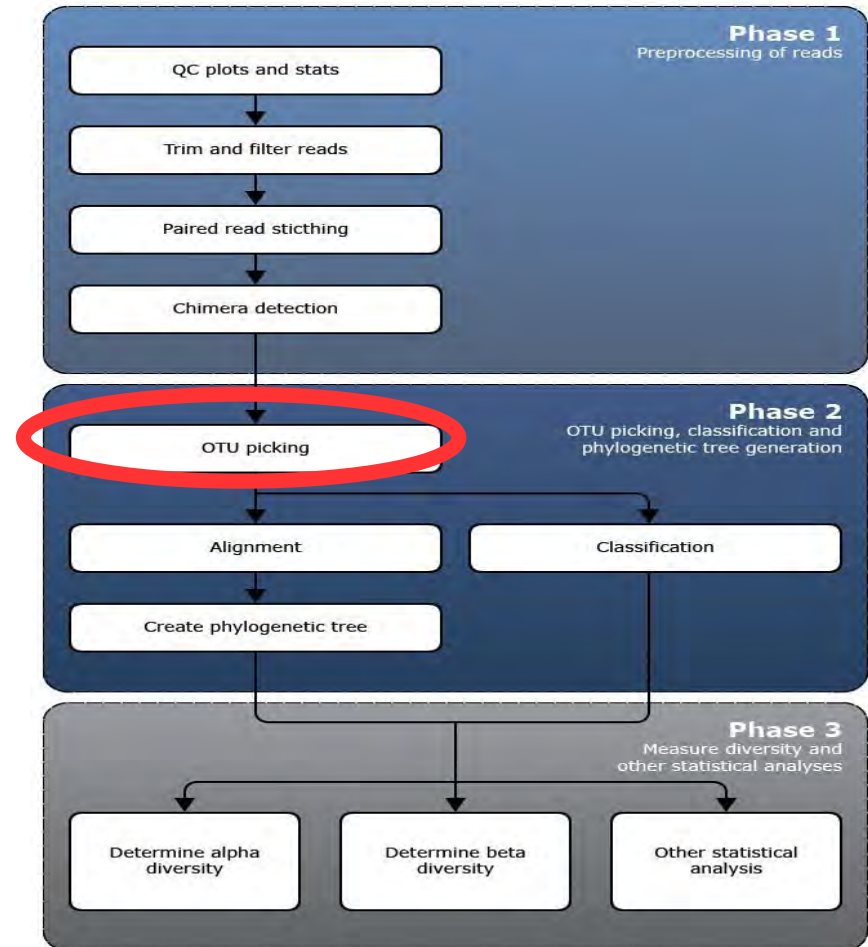
Pre-processing of input reads

- Demultiplexing
- Check for short read lengths
- Remove low quality bases
- Remove adapters and primers
- Remove chimeric sequences
- Merging of reads



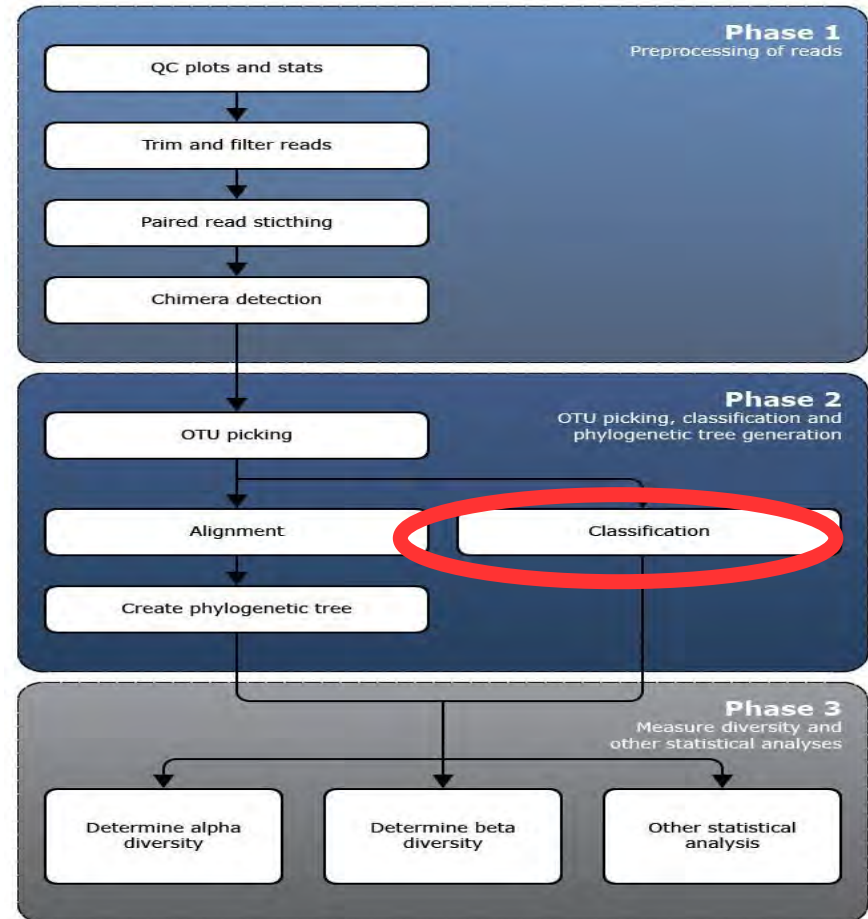
OTU picking

- An operational taxonomic unit is an operational definition of a species or group of species often used when only DNA sequence data is available.
- Clusters are formed based on sequence identity.
- 3 different approaches
 - De novo OTU picking
 - Closed reference OTU picking
 - Open reference picking
- A representative sequence is selected for downstream analysis.



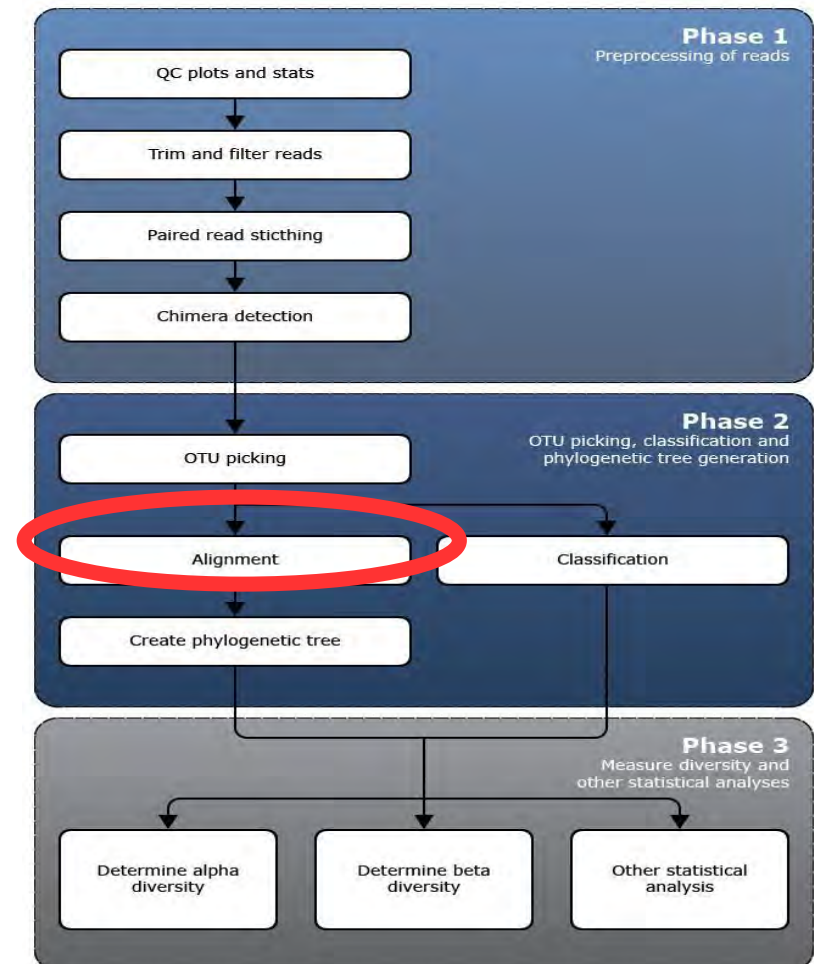
Classification

- A taxonomic identity is assigned to each representative sequence.
- Classification are done against three main reference databases with aligned, validated and annotated 16S rRNA genes: GreenGenes, Ribosomal Database Project (RDP) and Silva.



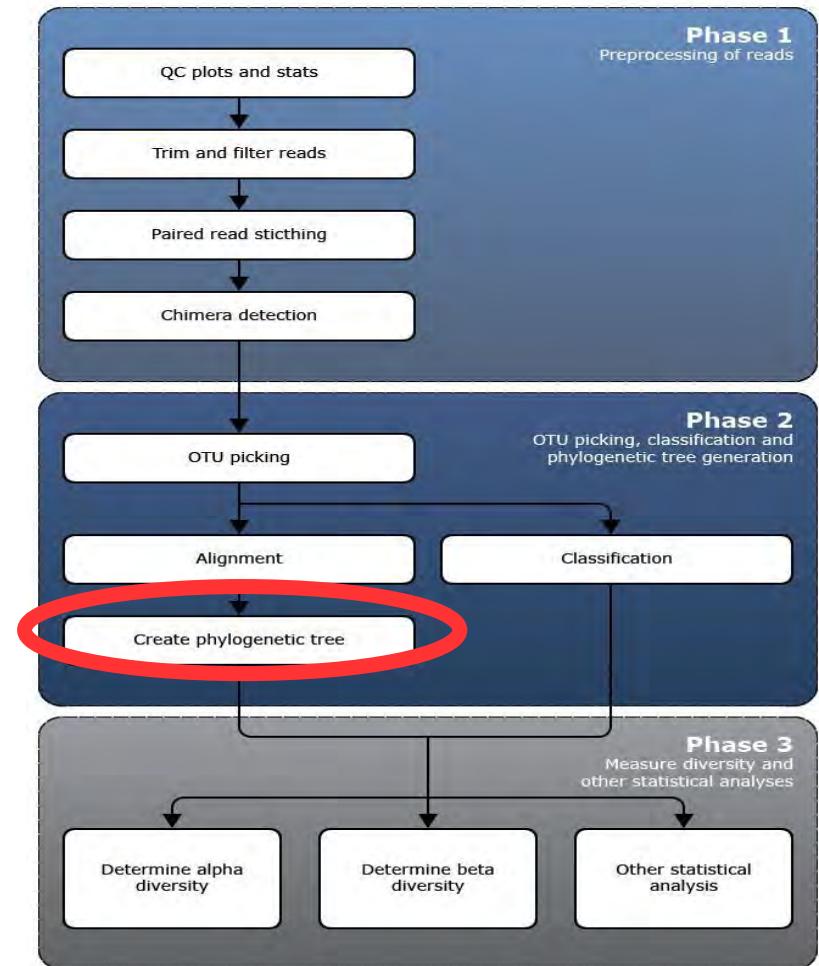
Alignment

- Alignment is the first step in generating a phylogenetic tree to understand the evolutionary relationship between samples.
- The aligners of choice are those that does alignments to a template (secondary structure) of the 16S sequence.



Create a phylogenetic tree

- The phylogenetic tree represents the relationship between the sequences in terms of the evolutionary distance from a common ancestor.
- In downstream analysis this tree is used for example in calculating the UniFrac distances and some alpha diversity measures.



Other 16S rRNA processing pipelines

4.3.1. Mothur Microbiome pipeline

APPLIED AND ENVIRONMENTAL MICROBIOLOGY, Dec. 2009, p. 7537–7541
0099-2240/09/\$12.00 doi:10.1128/AEM.01541-09
Copyright © 2009, American Society for Microbiology. All Rights Reserved.

Vol. 75, No. 23

Introducing mothur: Open-Source, Platform-Independent, Community-Supported Software for Describing and Comparing Microbial Communities[▽]

Patrick D. Schloss,^{1,2*} Sarah L. Westcott,^{1,2} Thomas Ryabin,¹ Justine R. Hall,³ Martin Hartmann,⁴
Emily B. Hollister,⁵ Ryan A. Lesniewski,⁶ Brian B. Oakley,⁷ Donovan H. Parks,⁸
Courtney J. Robinson,² Jason W. Sahl,⁹ Blaz Stres,¹⁰ Gerhard G. Thallinger,¹¹
David J. Van Horn,² and Carolyn F. Weber¹²

Department of Microbiology, University of Massachusetts, Amherst, Massachusetts¹; Department of Microbiology and Immunology, University of Michigan, Ann Arbor, Michigan²; Department of Biology, University of New Mexico, Albuquerque, New Mexico³; Department of Microbiology and Immunology, University of British Columbia, Vancouver, BC, Canada⁴; Department of Soil and Crop Sciences, Texas A&M University, College Station, Texas⁵; Department of Soil, Water, and Climate, University of Minnesota, St. Paul, Minnesota⁶; Department of Biological Sciences, University of Warwick, Coventry, United Kingdom⁷; Faculty of Computer Science, Dalhousie University, Halifax, NS, Canada⁸; Environmental Science and Engineering, Colorado School of Mines, Golden, Colorado⁹; Department of Animal Science, University of Ljubljana, Ljubljana, Slovenia¹⁰; Institute for Genomics and Bioinformatics, Graz University of Technology, Graz, Austria¹¹; and Department of Biological Sciences, Louisiana State University, Baton Rouge, Louisiana¹²

Received 30 June 2009/Accepted 26 September 2009

mothur aims to be a comprehensive software package that allows users to use a single piece of software to analyze community sequence data. It builds upon previous tools to provide a flexible and powerful software package for analyzing sequencing data. As a case study, we used mothur to trim, screen, and align sequences; calculate distances; assign sequences to operational taxonomic units; and describe the α and β diversity of

https://www.mothur.org/wiki/MiSeq_SOP

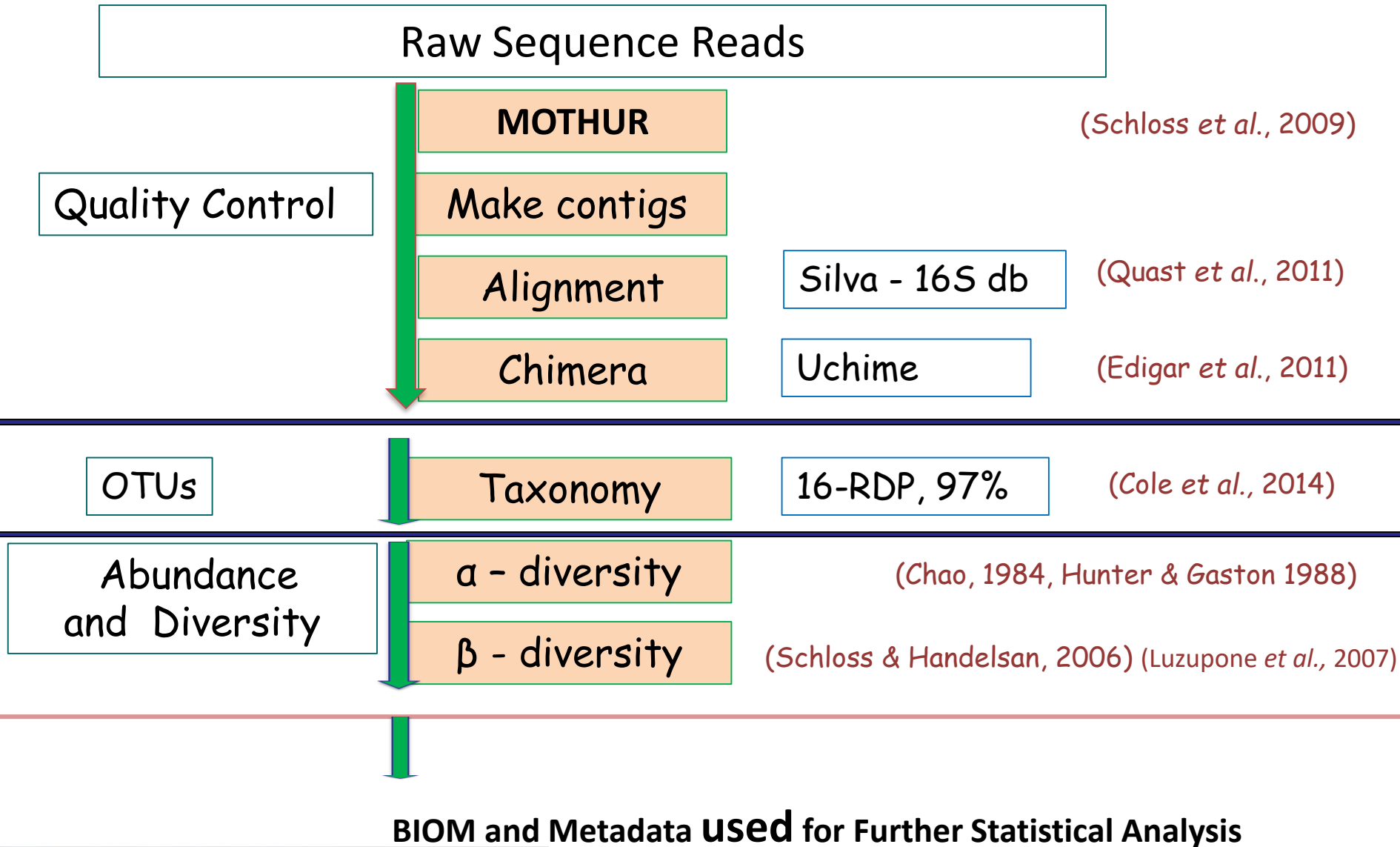


H3ABioNet

Pan African Bioinformatics Network for H3Africa

16SrRNA Intermediate Bioinformatics Online Course:
Int_BT_2019 Samson KM

MOTHUR - Analyses Workflow



QIIME PIPELINE



NIH Public Access

Author Manuscript

Not certified by peer review
Author manuscript; not certified by peer review

Published in final edited form as:

Nat Methods. 2010 May ; 7(5): 335–336. doi:10.1038/nmeth.f.303.

QIIME allows analysis of high-throughput community sequencing data

J Gregory Caporaso^{1,12}, Justin Kuczynski^{2,12}, Jesse Stombaugh^{1,12}, Kyle Bittinger³, Frederic D Bushman³, Elizabeth K Costello¹, Noah Fierer⁴, Antonio Gonzalez Peña⁵, Julia K Goodrich⁵, Jeffrey I Gordon⁶, Gavin A Huttley⁷, Scott T Kelley⁸, Dan Knights⁵, Jeremy E Koenig⁹, Ruth E Ley⁹, Catherine A Lozupone¹, Daniel McDonald¹, Brian D Muegge⁶, Meg Pirrung¹, Jens Reeder¹, Joel R Sevinsky¹⁰, Peter J Turnbaugh⁶, William A Walters², Jeremy Widmann¹, Tanya Yatsunenko⁶, Jesse Zaneveld², and Rob Knight^{1,11}

Rob Knight: rob.knight@colorado.edu

Developed by Caporaso et al., 2010; <https://QIIME2.org/>



H3ABioNet

Pan African Bioinformatics Network for H3Africa

16SrRNA Intermediate Bioinformatics
Online Course: Int_BT_2019 Samson KM

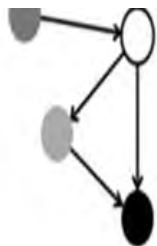
What is QIIME?

Quantitative Insight Into Microbial Ecology

- It is a Next-generation Microbiome Bioinformatics Platform that is **Extensible, Free, Open source and Community Developed**.
- It is a powerful suite of microbiome analysis packages where researchers can start an analysis with **raw DNA sequence data and finish with publication-quality figures and statistical results**
- **QIIME2** is a complete **redesign** of the QIIME1 while retaining the features that makes it a powerful and widely-used tool

QIIME2 Key Features

Tracking and Exploration

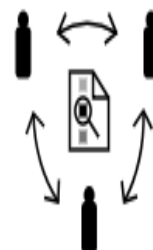


Automatically track your analyses with decentralized data provenance — no more guesswork on what commands were run!



Interactively explore your data with beautiful visualizations that provide new perspectives.

Easy sharing, Suite of tools



Easily share results with your team, even those members without QIIME 2 installed.



Plugin-based system — your favorite microbiome methods all in one place.

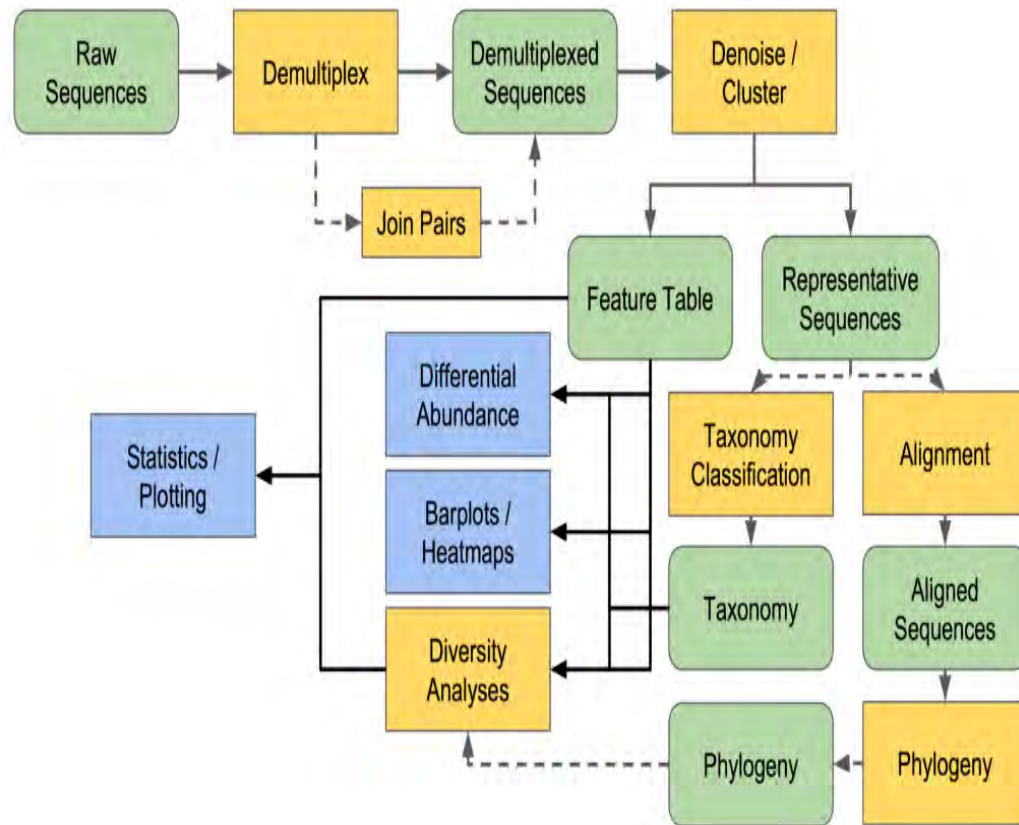


H3ABioNet

Pan African Bioinformatics Network for H3Africa

Module 4: Part IV.. cont..... Bioinformatics Pipeline

QIIME2 CONCEPTUAL WORKFLOW

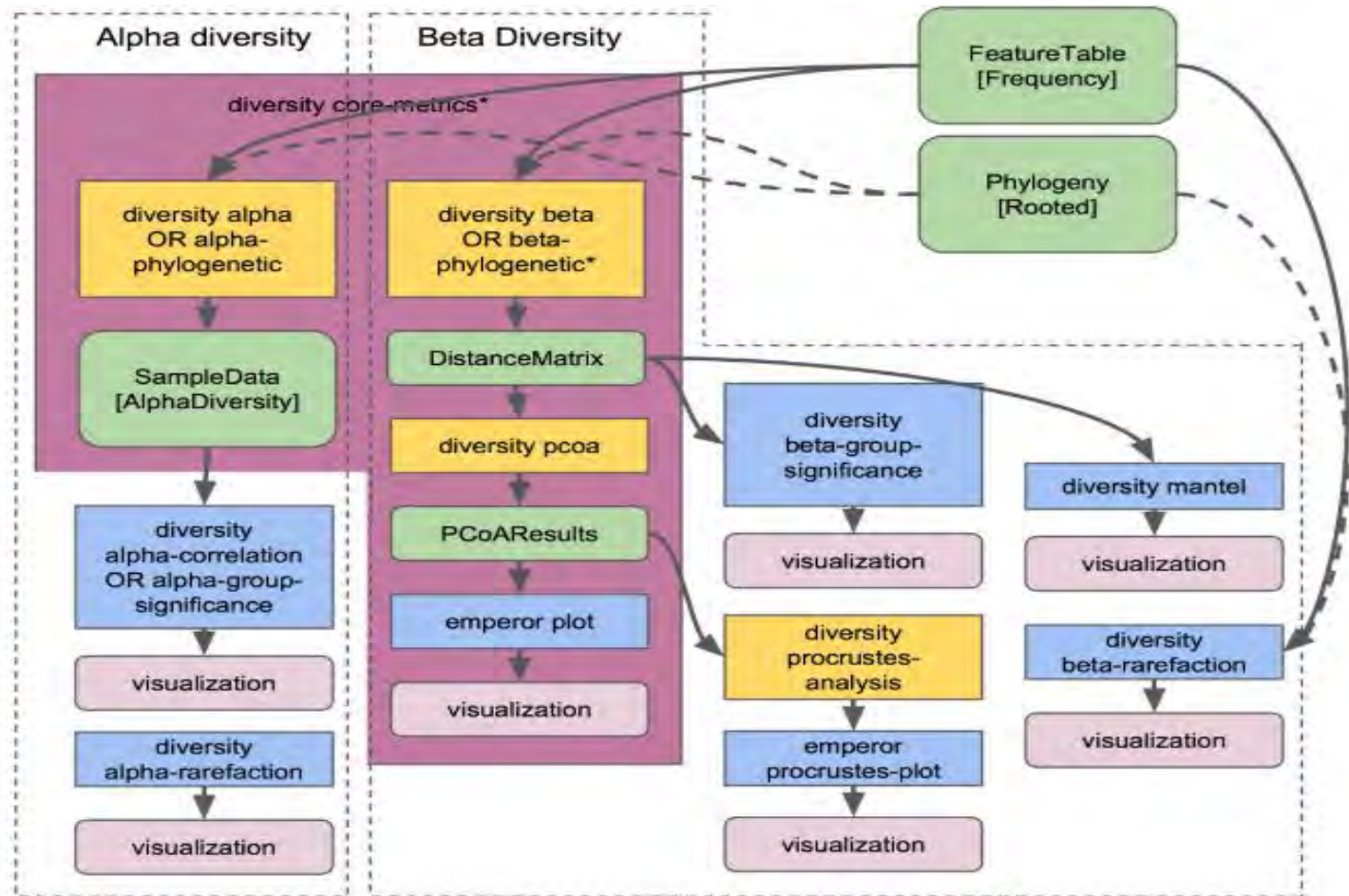


Data are imported as a QIIME2 artifact to be used by a QIIME2 action (except the metadata).

Users may enter the workflow at different stages. Most will have raw sequence (e.g., FASTQ or FASTA) data, which should be imported.

Other users may start with demultiplexed sequence data or even a feature table

Alpha and Beta Diversity Analysis



q2-diversity plugin.....



H3ABioNet

Pan African Bioinformatics Network for H3Africa

Example of Selected q2-plugins

Analyse longitudinal microbiome data:

q2-longitudinal plugins: Performs statistical analyses of longitudinal studies,

- ◇ i.e., where samples are collected from patients/subjects/sites repeatedly over time.

Predict the future (or the past)

q2-sample-classifier plugin: predict sample metadata as a function of feature data

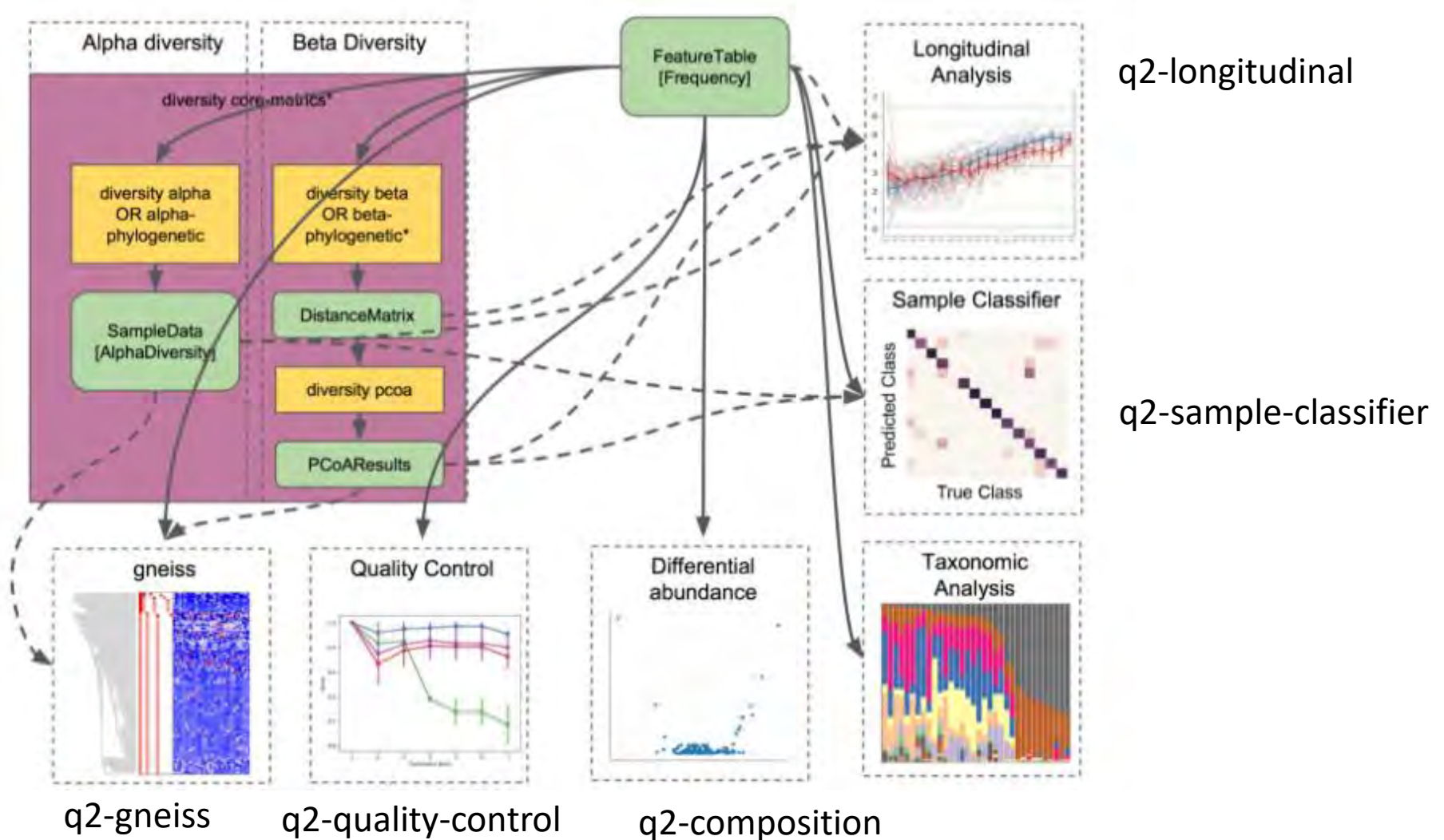
- ◇ i.e. can we use a fecal sample to predict cancer susceptibility? Or predict wine quality based on the microbial composition of grapes before fermentation?

Differential abundance

q2-composition plugin (based on ANCOM) and **q2-gneiss** plugin

- ◇ Determines which features are significantly more/less in different groups of samples

QIIME2 Statistical Analysis Types



Example (1) Plots- Microbial Profile



Download

SVG (bars) SVG (legend) CSV

Taxonomic Level

Level 3

Color Palette

schemeAccent

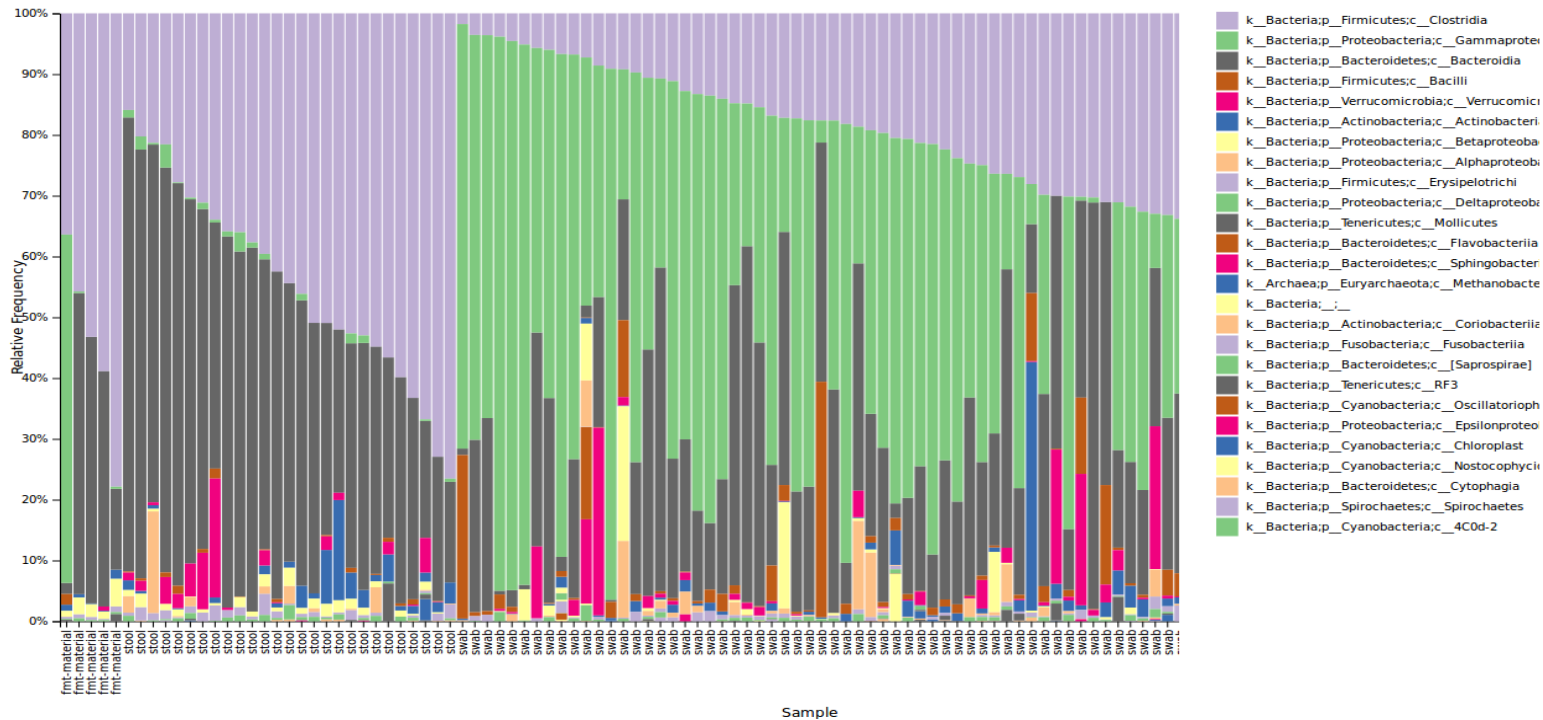
Sort Samples By

sample-type

Ascending

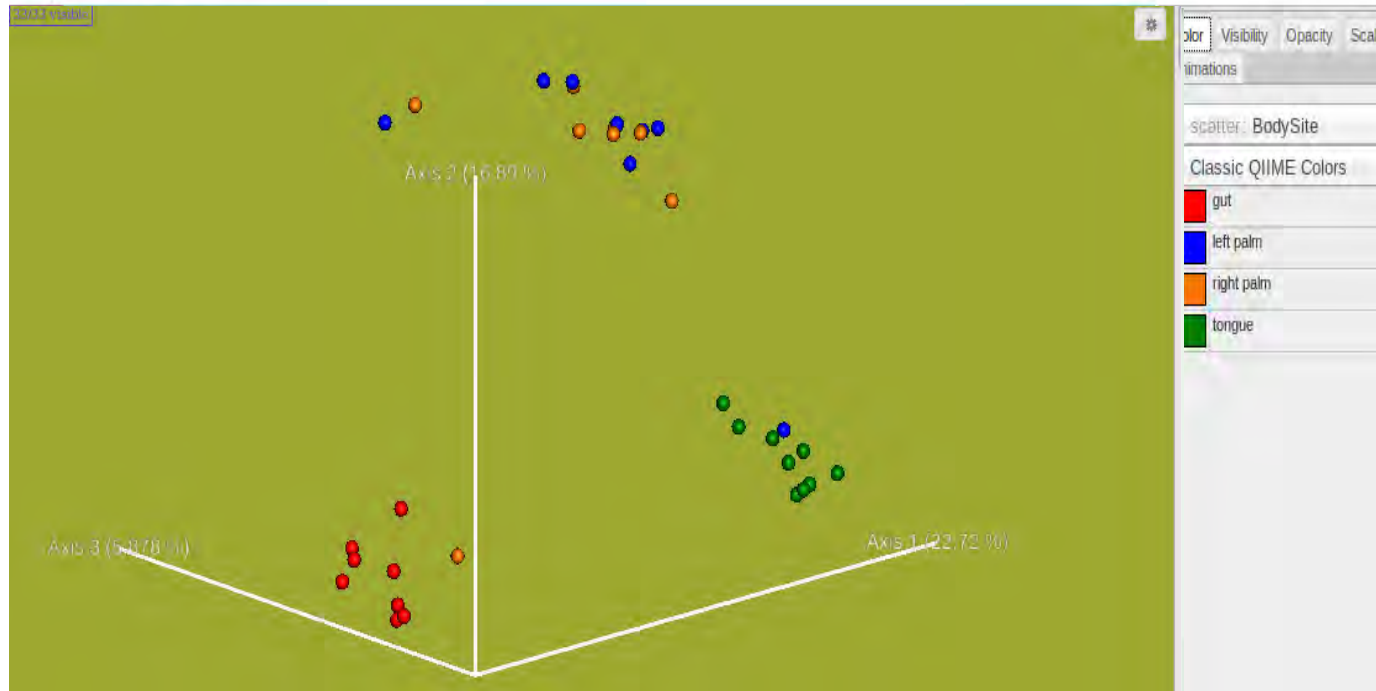
Relabel X? ☒

swab | k__Bacteria;p__Proteobacteria;c__Gammaproteobacteria | 46.645%



Example (2) Plots- beta diversity - PCoA plots

qiime2



QIIME2 cont...

Further readings

<https://qiime2.org/>

Installation of qiime2 here

<https://docs.qiime2.org/2019.7/install/>

<https://docs.qiime2.org/2019.7/>

Variety of tutorials found here...

<https://docs.qiime2.org/2019.7/tutorials/>

Subscribe to qiime2 forum for assistance here..

<https://forum.qiime2.org/>

DADA2 16S rRNA Bioinformatic Pipeline

Nature methods

Author Manuscript

HHMI Public Access

DADA2: High resolution sample inference from Illumina amplicon data

Benjamin J Callahan, Paul J McMurdie, [...], and Susan P Holmes

[Additional article information](#)

Associated Data

- [Supplementary Materials](#)

Abstract

We present DADA2, a software package that models and corrects Illumina-sequenced amplicon errors. DADA2 infers sample sequences exactly, without coarse-graining into OTUs, and resolves differences of as little as one nucleotide. In

Developed by Callahan et al., 2016



H3ABioNet

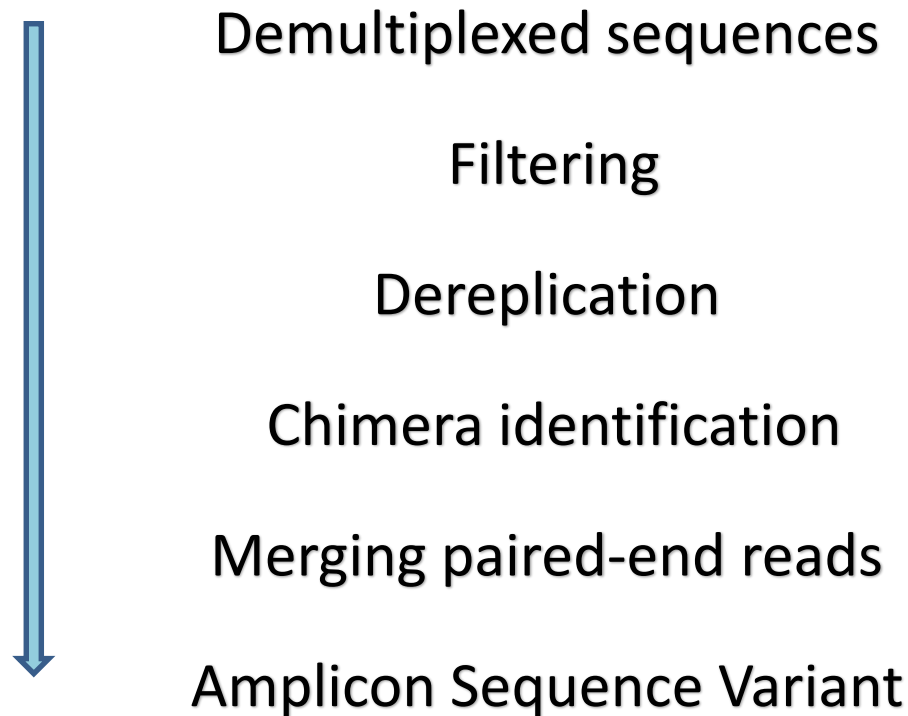
Pan African Bioinformatics Network for H3Africa

DADA2 cont...

- **DADA2** stands for - **D**ivisive **A**mplicon **D**enoising **A**lgorithm
- It is an **Open source** R-based package which corrects amplicons errors directly without constructing OTUs
- It infers sample sequences exactly and **resolves differences of as little as one nucleotide**
- The output are the Sequence Amplicon Variant (ASV) also named Exact Sequence Variant (ESV), as opposed to the OTUs
- In vie of the potential of this pipeline, the **H3ABioNet** **have developed DADA2 Nextflow Pipeline** for the Microbiome data processing. You will use this pipeline in practical session.

In view of that, The H3ABioNet have a developed a DADA2 Nextflow pipeline

The DADA2 workflow



Import into R Phyloseq package for downstream Analysis



DADA2 cont...

Further reading on DADA2

Installation of DADA2

<https://bioconductor.org/packages/devel/bioc/vignettes/dada2/inst/doc/dada2-intro.html>

Tutorials available here:

<https://benjjneb.github.io/dada2/tutorial.html>

Module 4: 16S Microbiome Pipeline Summary

16S rRNA diversity vs shotgun analysis

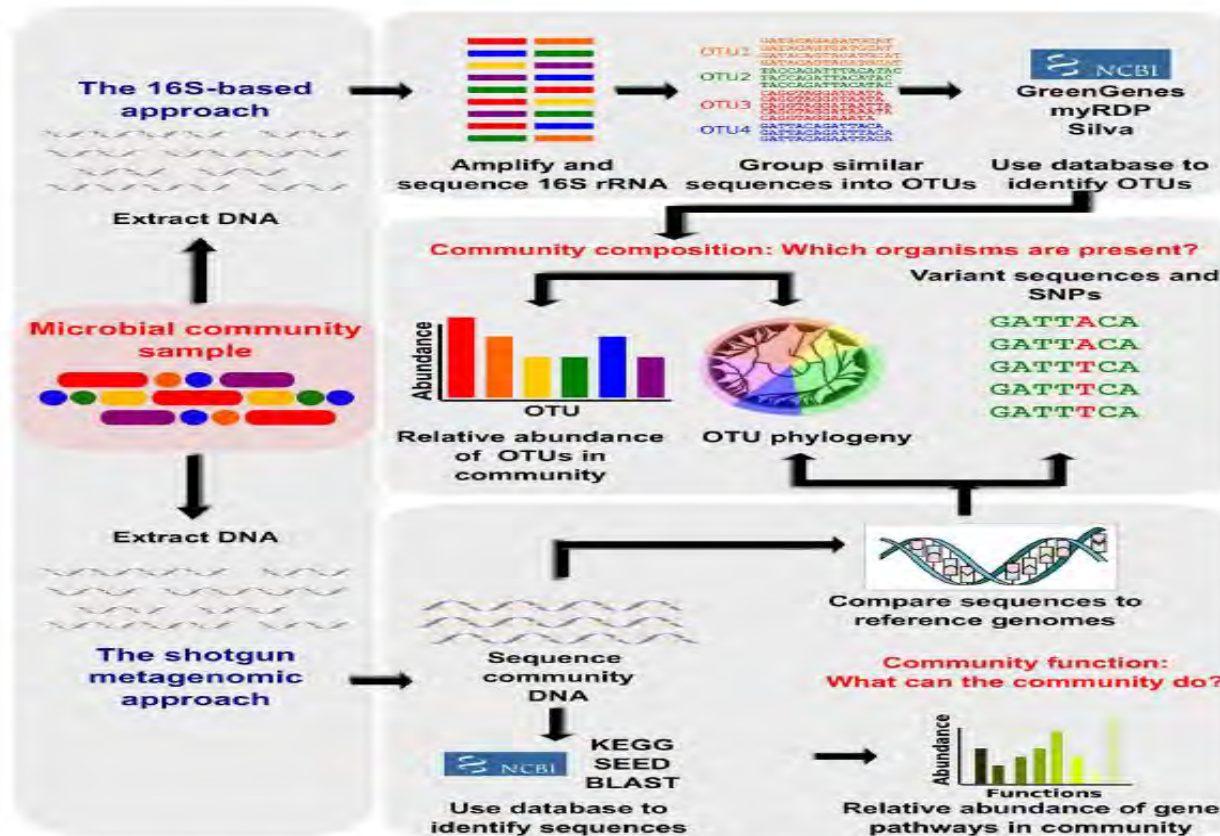
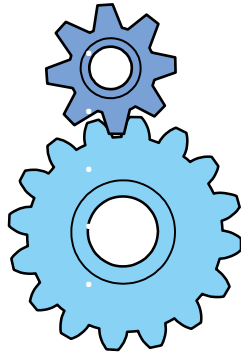


Figure 1. Bioinformatic methods for functional metagenomics. Studies that aim to define the composition and function of uncultured microbial communities are often referred to collectively as “metagenomic,” although this refers more specifically to particular sequencing-based assays. First, community DNA is extracted from a sample, typically uncultured, containing multiple microbial members. The bacterial taxa present in

SUMMARY

Raw sequence data

16S rRNA PIPELINE



BIOM file
otus_repsetOUT.fa
Phylogenetic tree
Taxonomy file
otu_table

STATISTICIAN



H3ABioNet

Pan African Bioinformatics Network for H3Africa

Other available pipelines

- UPARSE:
<http://www.drive5.com/uparse>
- IM Tornado: <https://github.com/pjeraldo/imtornado2>
- FROGS: <https://github.com/geraldinepascal/FROGS>
- VSEARCH: <https://github.com/torognes/vsearch>

Acknowledgements

H3ABioNet for some contents of these slides

The CBIO at UCT for the inhouse microbiome-pipeline

Verena Ras and Gerrit Botha, for their devoted time
on this training

Dar es Salaam Institute of Technology : My employer

Thank you

Asanteni