# Intro to Amplicon Sequence Variants Analysis

Ashley Shade

EDAMAME 29 June 2018

# The "OTU"
## operational taxonomic unit

- Species = basic unit of classification
- Defined somewhat arbitrarily
- Typical = 97% sequence identity
  - Originally, identity based on *full length* 16S rRNA gene
  - roughly equivalent to genus level
  - Does not well-distinguish "taxa" for all bacteria (*e.g., Streptomyces*)
  - Used in part to minimize species inflation by including amplification and sequencing errors
- Different methods of defining OTUs will result in different numbers of taxa! Different numbers of taxa mean different perspectives of diversity!

# Approaches to Picking OTUs

- Reference based :  percent identity to defined taxa populating a reference database
  - Pros:  you know the taxa are "real"!
  - Cons:  Weird environments don't have many representatives in databases, only as good as your database, could end up throwing out a lot of real stuff, use of "representative sequences" of the cluster

- De novo :  percent identity to other sequences in the dataset; taxonomic assignment to the OTUs happens afterwards
  - Pros:  Good for weird environments with low representation in databases
  - Cons:  Computationally expensive, "greedy" algorithms can artificially inflate diversity

- Open reference : cluster against a reference db first, and anything that doesn't hit gets clustered de novo
  - Best of both worlds?  Now can optimized so that new de novo OTUs are added to the original database and used subsequently in "reference" clustering
  - See Rideout et al. 2014 PeerJ

# What are amplicon sequence variants?

- Essentially, OTUs clustered at 100% sequence identity, after *rigorous and conservative* quality control of the data

- Direct consequences:
  - OTU inflation ?
    - Multiple 16S rRNA copies from the same organism may get split into different OTUs
    - If the quality control is too relaxed and keeps trashy data
  - OTU deflation?
    - If the quality control algorithm is too conservative and throws out real data
    - If there were a LOT of errors in the data and they were "corrected"

# Other names for sequence variants

- Exact sequence variants (ESV) - Glassman and Martiny 2018
- Amplicon sequence variants (ASV) – Callahan and Holmes / Phyloseq folks
- Sequence variants (SV) –- QIIME folks
- ZOTUs (zero-radius OTUs) – usearch/unoise, Edgar et al.
- Oligotypes – Eren et al.

Robust to ecological patterns

CSH Cold Spring Harbor Laboratory

# bioRxiv
THE PREPRINT SERVER FOR BIOLOGY

New Results

## Ecological patterns are robust to use of exact sequence variants versus operational taxonomic units

Sydney I Glassman, Jennifer BH Martiny

Now accepted for publication in *mSphere*

Abstract    Info/History    Metrics                    Preview PDF

## Abstract

Recent controversy focuses on the best method for delineating microbial taxa, based on either traditional operational taxonomic units (OTUs) or exact sequence variants (ESVs) of marker gene sequences. We sought to test if the binning approach (ESVs versus OTUs defined by >97% sequence similarity) affected the conclusions of a large field study. The dataset included sequences targeting all bacteria (16S) and fungi (ITS), across multiple environments diverging

Work suggesting that "tighter" clusters should be used instead of 97%

# Advantages of using sequence variants

- Can be re-used and compared directly across studies: "consistent labels"

- For some taxa, represent the most taxonomically resolved autonomous unit

- Some quality control and "clustering" methods afford a high level of confidence
  - Note that the most rigorous of these methods are not really clustering in the traditional sense

- "de novo", so not reliant on limited databases, but precisely linked to any reference at 100%

# Disadvantages of ASV

- Quality control **IS SO IMPORTANT** – errors can propagate within and across datasets quickly if not careful about QC

- For some taxa, other less resolved clustering (e.g., 97%) is biologically appropriate for ecologically meaningful taxa

- For some taxa, (e.g., Streptomyces), ASV are no better at resolving than 97%

- The most rigorous methods of ASV analysis are more computationally expensive

- It is cheap and easy to re-pick OTUs/ASVs for each new study/meta-analysis, and arguably needs to be done anyway to ensure integrity and rigor

- Sequencing effort matters, and must be consistent across studies

# Denoising: essential for ASV
# Errors in amplicon HTS

- Amplification errors
    - Substitutions – incorrect base pairings (*point error*)
    - Gaps – polymerase slips (*point error*)
    - Chimeras – combined biological templates by polymerase

- Sequencing errors
    - Insertions and deletions ("indels", *point errors*)
    - Spurious base calls (*point error*)

- Contaminants
    - In reagents or kits (or from researcher!)
    - "Cross-talk" between Illumina lanes

Edgar 2018 bioRXiv

# Tools for determining sequence variants

- **Oligotyping/ MEN** (Eren et al. 2013) – includes quality filtering but works best when pre-QC steps are taken

- **QIIME2**: DEBLUR for QC using information theory method, and then q2-vsearch for 100% clustering

- **unoise3**: "zOTUs" – includes quality filtering, pools all samples to error correct using heuristic model and determine variants

- **DADA2**: "amplicon sequence variants" - includes quality filtering; uses single samples to build an error correction model

# Comparing the denoising methods

## Denoising the Denoisers: An independent evaluation of microbiome sequence error-correction methods

Research article   Genomics   Microbiology   Data Science

Jacob T Nearing[1], Gavin M Douglas[1], André M Comeau[2], Morgan G.I Langille[1,3]

February 23, 2018

Correcting
for
ribosomal
copy
number?

## Correcting for 16S rRNA gene copy numbers in microbiome surveys remains an unsolved problem

Stilianos Louca,[1,2] Michael Doebeli,[1,2,3] and Laura Wegener Parfrey[1,2,4]

Author information ▸ Article notes ▸ Copyright and License information ▸ Disclaimer

### Abstract

Go to: ⊙

The 16S ribosomal RNA gene is the most widely used marker gene in microbial ecology. Counts of 16S sequence variants, often in PCR amplicons, are used to estimate proportions of bacterial and archaeal taxa in microbial communities. Because different organisms contain different 16S gene copy numbers (GCNs), sequence variant counts are biased towards clades with greater GCNs. Several tools have recently been developed for predicting GCNs using phylogenetic methods and based on sequenced genomes, in order to correct for these biases. However, the accuracy of those predictions has not been independently assessed. Here, we systematically evaluate the predictability of 16S GCNs across bacterial and archaeal clades, based on ~ 6,800 public sequenced genomes and using several phylogenetic methods. Further, we assess the accuracy of GCNs predicted by three recently published tools (PICRUSt, CopyRighter, and PAPRICA) over a wide range of taxa and for 635 microbial communities from varied environments. We find that regardless of the phylogenetic method tested, 16S GCNs could only be accurately predicted for a limited fraction of taxa, namely taxa with closely to moderately related representatives (≤15% divergence in the 16S rRNA
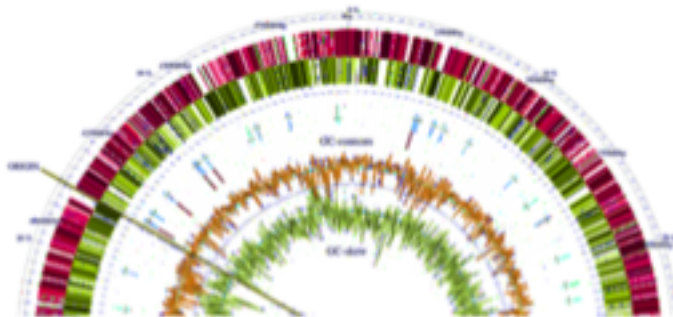
# Tools are being curated!

## Oligotypes: the classic approach to ASV

## Oligotyping: differentiating between closely related microbial taxa using 16S rRNA gene data

A. Murat Eren*, Lois Maignien, Woo Jun Sul, Leslie G. Murphy, Sharon L. Grim, Hilary G. Morrison and Mitchell L. Sogin

*Josephine Bay Paul Center for Comparative Molecular Biology and Evolution, Marine Biological Laboratory, Woods Hole, MA, 02543, USA*

### Summary

1. Bacteria comprise the most diverse domain of life on Earth, where they occupy nearly every possible ecological niche and play key roles in biological and chemical processes. Studying the composition and ecology of bacterial ecosystems and understanding their function are of prime importance. High-throughput sequencing technologies enable nearly comprehensive descriptions of bacterial diversity through 16S ribosomal RNA gene amplicons. Analyses of these communities generally rely upon taxonomic assignments through reference data bases or clustering approaches using *de facto* sequence similarity thresholds to identify operational taxonomic units. However, these methods often fail to resolve ecologically meaningful differences between closely related organisms in complex microbial data sets.

2. In this paper, we describe oligotyping, a novel supervised computational method that allows researchers to investigate the diversity of closely related but distinct bacterial organisms in final operational taxonomic units identified in environmental data sets through 16S ribosomal RNA gene data by the canonical approaches.

3. Our analysis of two data sets from two different environments demonstrates the capacity of oligotyping at discriminating distinct microbial populations of ecological importance.

# Oligotyping

- Capitalizes on most "informative" variable sites along a nucleotide sequence
  - "systematic" identification of nucleotide positions that are informative
- Premise
  - Distinguish sequencing errors that result in insertions or deletions (indels) from real mutations that result in sequence variation from a shared evolutionary history
  - Key component:  reads must be the same length (good for Illumina) and, if possible **high quality**
  - Uses Shannon Entropy to determine informative sequence variants

# Shannon Entropy

- Quantify the uncertainty at each nucleotide position
- Think about: the amount of information that can be gleaned from a nucleotide position.
  - Information has an element of "surprise"
- Probability distribution: one event has a probability of 1, and all else of 0.
- Maximum entropy = log(n), where n is number of events
  - Events are like surprises – shifts in nucleotide content at the same nucleotide position, across similar sequences
  - For a nucleotide position that contains all 4 bases, the entropy is log(4)
  - For a position that contains only 2 bases, the entropy is log(2)
  - For a position that contains 1 base, the entropy is log(1) == 0
- Positions with maximum entropy are flagged as informative for oligotyping

Oligotyping workflow

# Oligotyping is a supervised, iterative process



Fig. 1. Flowchart of an oligotyping analysis

http://merenlab.org/software/oligotyping/

# ORIGINAL ARTICLE

# Minimum entropy decomposition: Unsupervised oligotyping for sensitive partitioning of high-throughput marker gene sequences

A Murat Eren, Hilary G Morrison, Pamela J Lescault, Julie Reveillaud, Joseph H Vineis and Mitchell L Sogin

*Josephine Bay Paul Center for Comparative Molecular Biology and Evolution, Marine Biological Laboratory, Woods Hole, MA, USA*

Molecular microbial ecology investigations often employ large marker gene datasets, for example, ribosomal RNAs, to represent the occurrence of single-cell genomes in microbial communities. Massively parallel DNA sequencing technologies enable extensive surveys of marker gene libraries that sometimes include nearly identical sequences. Computational approaches that rely on pairwise sequence alignments for similarity assessment and *de novo* clustering with *de facto* similarity thresholds to partition high-throughput sequencing datasets constrain fine-scale resolution descriptions of microbial communities. Minimum Entropy Decomposition (MED) provides a computationally efficient means to partition marker gene datasets into 'MED nodes', which represent homogeneous operational taxonomic units. By employing Shannon entropy, MED uses