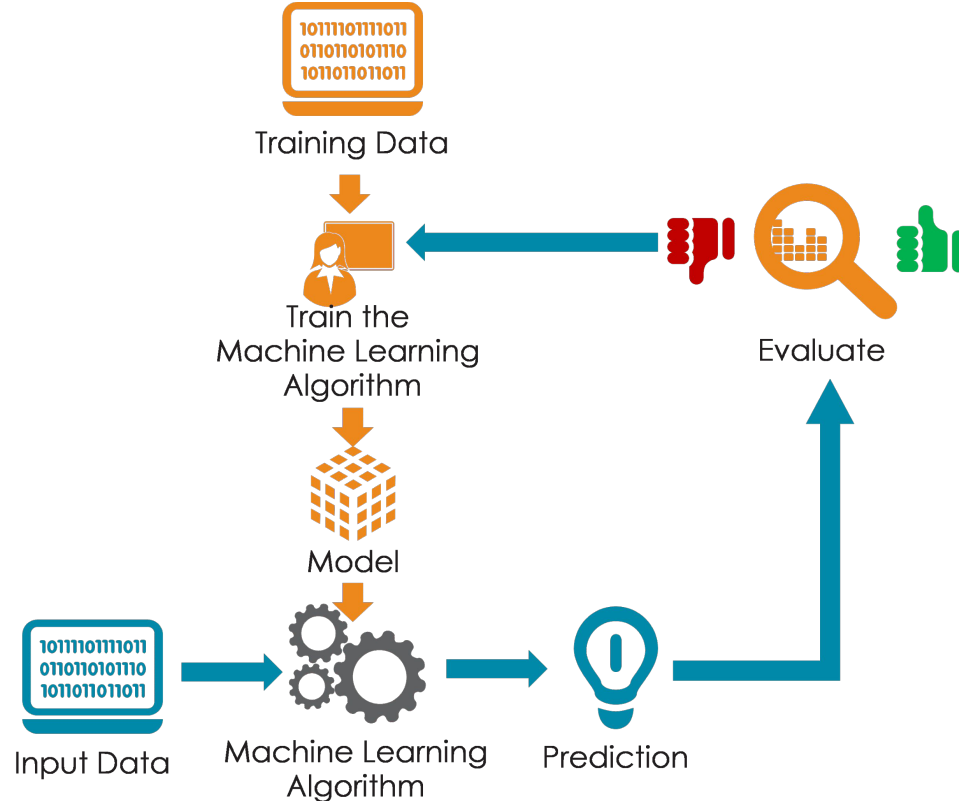


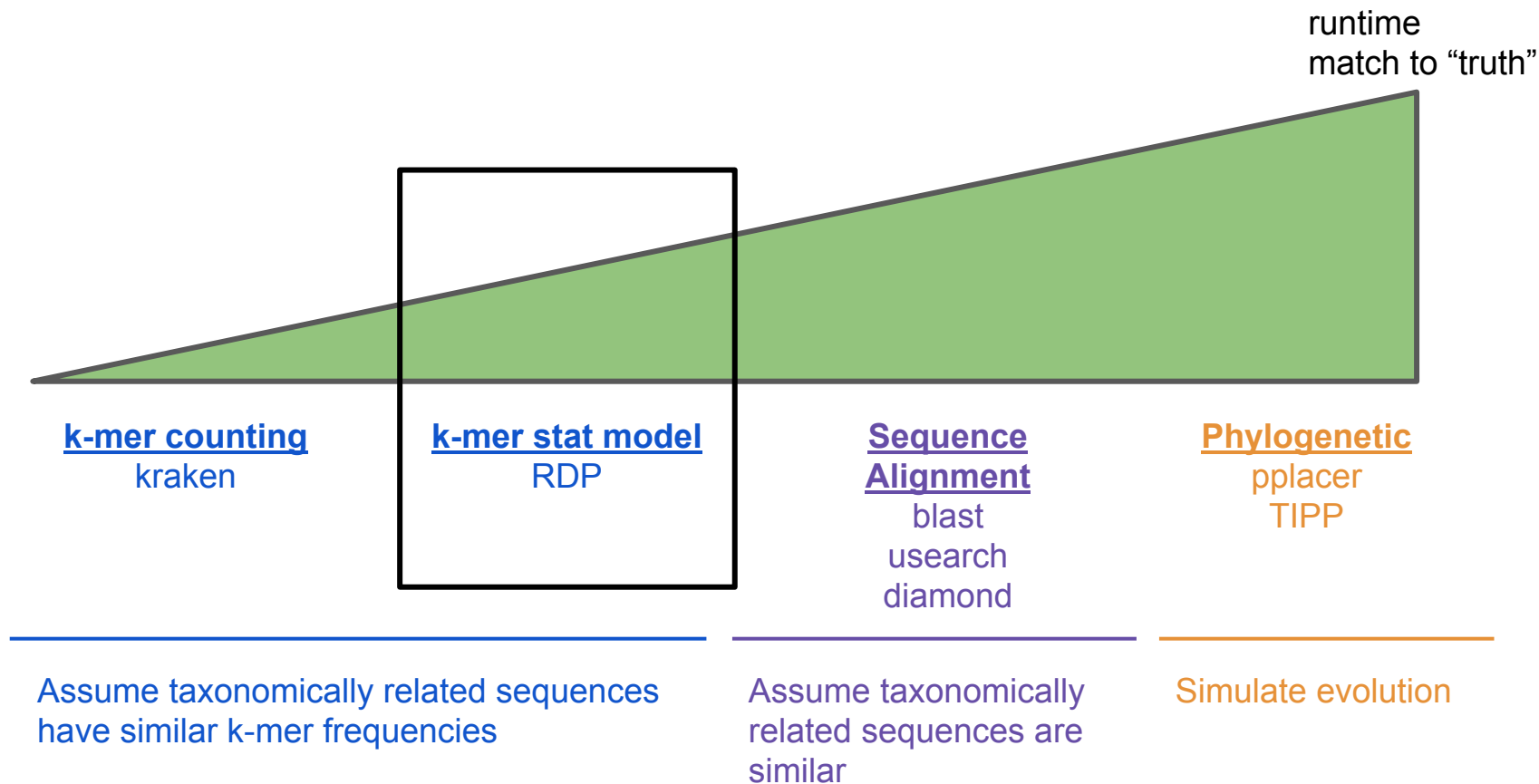
Machine Learning Approaches: Ribosomal Database Classifier

M³ Workshop
January 10, 2019

Using machine learning to predict taxonomic labels



Landscape of Taxonomic Approaches



Ribosomal Database Project Classifier



- <https://rdp.cme.msu.edu/>
- Github: <https://github.com/rdpstaff/classifier>
- Publication: [Naive Bayesian classifier for rapid assignment of rRNA sequences into the new bacterial taxonomy. Wang et al. Appl Environ Microbiol 2007.](#)
- Default for mothur's `classify.seqs` (method=wang) and option for QIIME's `assign_taxonomy.py`

Building the RDP Classifier

- Uses **k-mers**-- all possible substrings of length k that are contained in a string

Sequence: **ATGGAAGTCGCGGAA**

8-mers
ATGGAAGT
TGGAAGTC
GGAAGTCG
GAAGTCGC
AAGTCGCG
AGTCGCGG
GTCGCGGA
TCGCGGAA

Building the RDP Classifier

- Uses **k-mers**-- all possible substrings of length k that are contained in a string

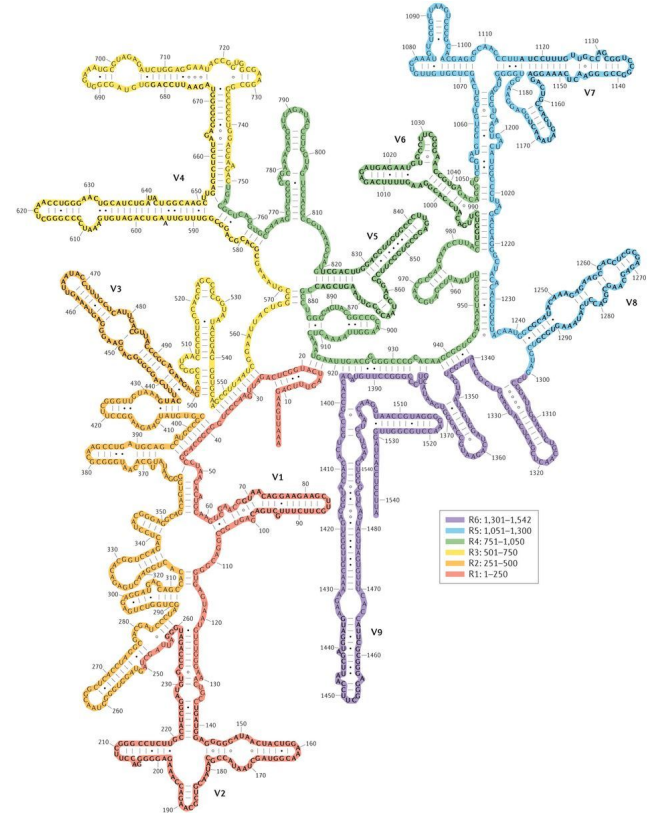
Sequence: **ATGGAAGTCGCGGAA**

8-mers
ATGGAAGT
TGGAAGTC
GGAAGTCG
GAAGTCGC
AAGTCGCG
AGTCGCGG
GTCGCGGA
TCGCGGAA

- Identifies all possible 8-mers in a taxon database and builds a table of how many times each 8-mer appears in each taxon

RDP taxonomy database is based on marker genes

- Why marker genes?
 - Compare apples to apples
 - Avoid lateral gene transfer
 - Can (in theory) extrapolate based on evolutionary principles
 - Can target experimentally (\$'s matter)
- RDP has classifiers for:
 - 16S ribosomal RNA
 - Fungal LSU
 - Fungal ITS
- You can also train your own classifier with updated taxonomy or different marker genes!



Using the RDP Classifier to assign taxonomy

- Calculates all 8-mers in a query sequence and looks them up in the reference table to calculate the probability that the query sequence is a member of a genus

By Bayes' theorem, the probability that an unknown query sequence, S , is a member of genus G is

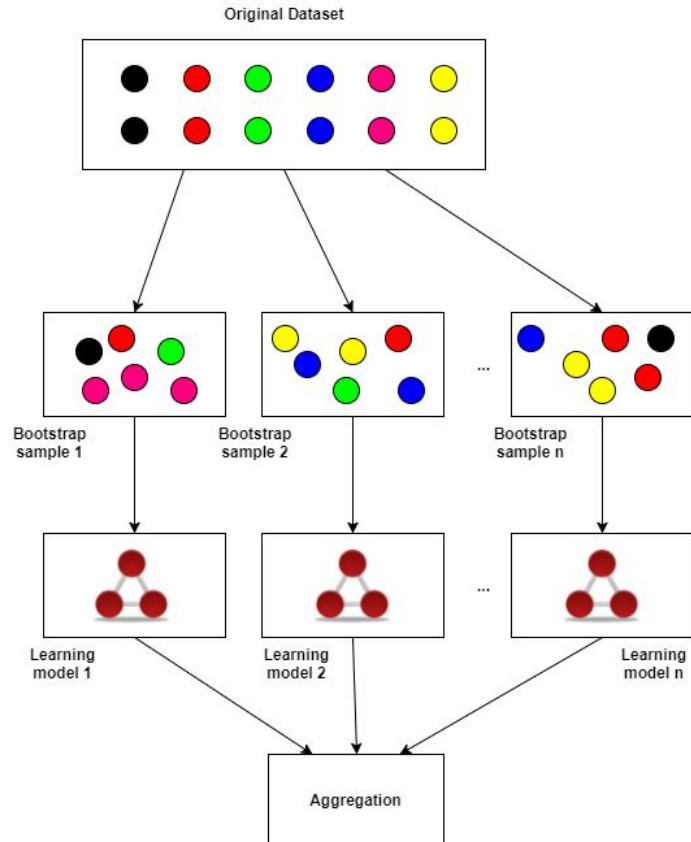
$$P(G|S) = P(S|G) \times P(G)/P(S)$$

$P(G)$ = the prior probability of a sequence being a member of G

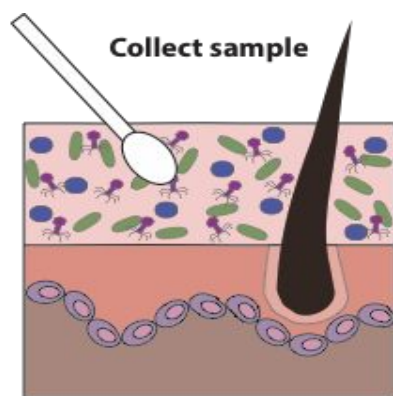
$P(S)$ = the overall probability of observing sequence S (from any genus)

$P(S|G)$ = the joint probability of observing from genus G a (partial) sequence, S , containing a set of words, $V = \{v_1, v_2, \dots, v_f\} = \prod P(v_i|G)$

Assigning confidence to taxonomy predictions



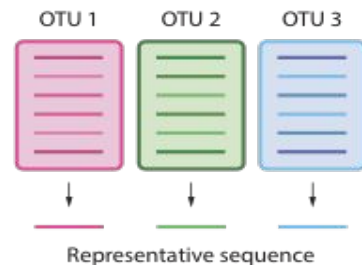
- Use random-resampling (bootstrapping) to assign confidence
- Runs 100 times; each time it chooses a subset of $\frac{1}{8}$ of all possible 8-mers from the query sequence
- For a particular taxonomic assignment, confidence is assessed by the number of times a particular genus is assigned out of 100 trials
- Confidence score ranges from 0 (not confident) - 1 (very confident)
- User can set a minimum confidence score [default: 0.8]



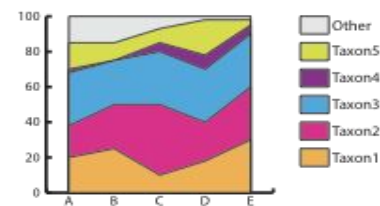
**PCR Amplify and Sequence
16S rRNA Marker Gene**



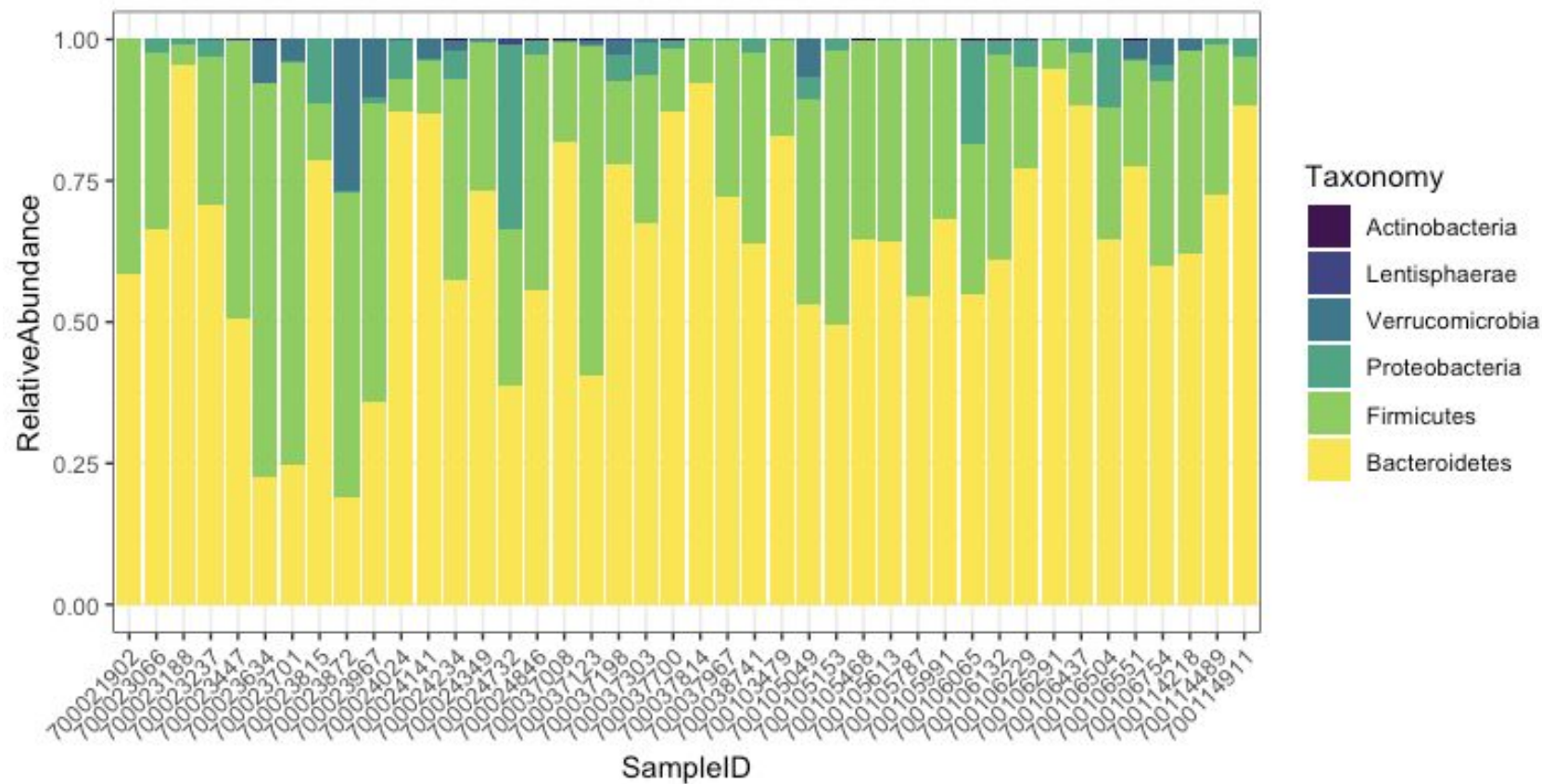
**Group similar sequences into
Operational Taxonomic Units (OTUs)**



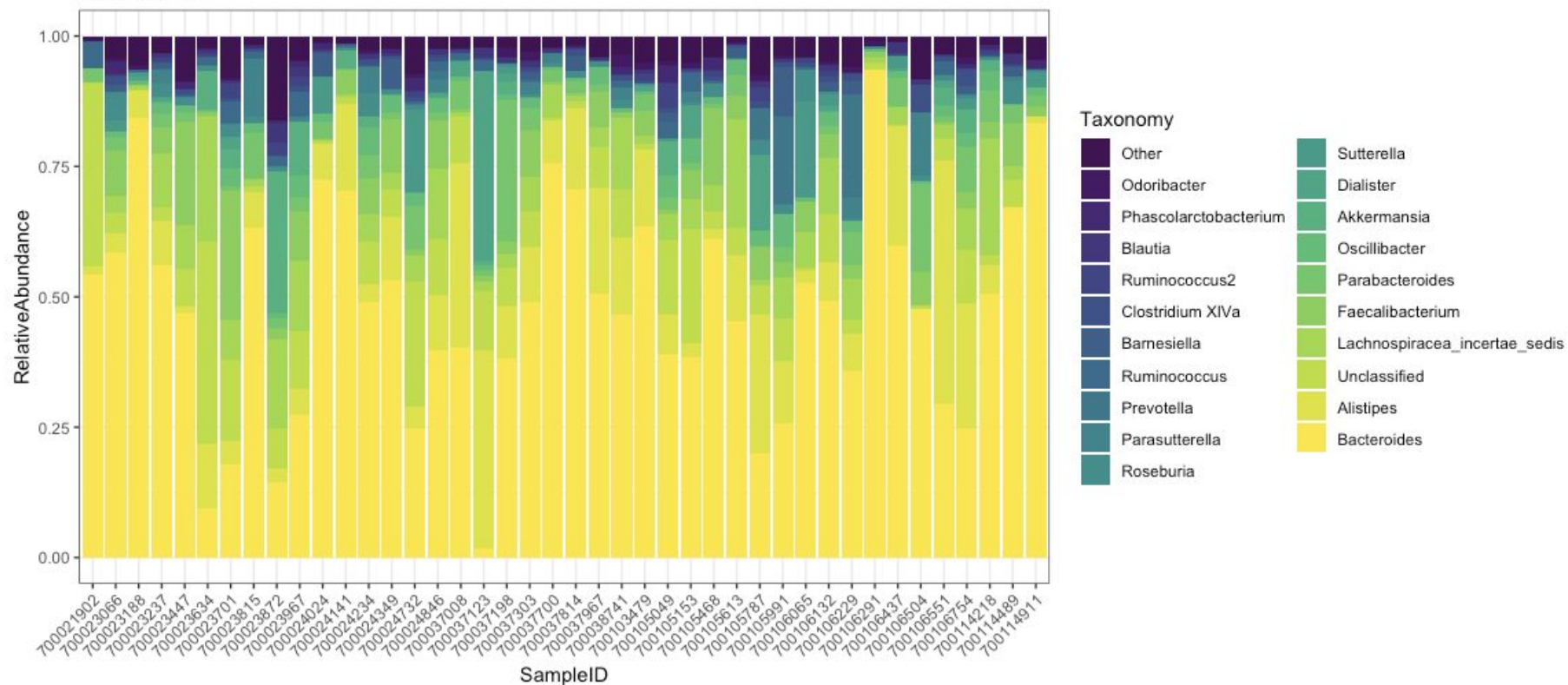
**Compare to
Reference Databases**



RDP-Phylum



RDP-Genus



If you were to train a classifier, what features would you use?