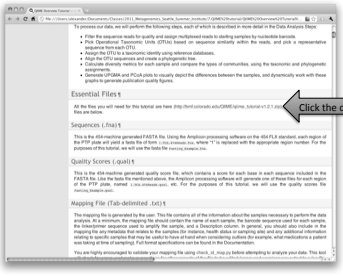


## QIIME tutorial

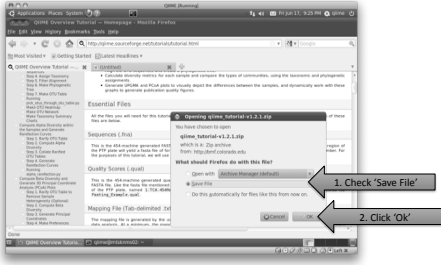
<http://qiime.org/tutorials/tutorial.html>

## Download the necessary files



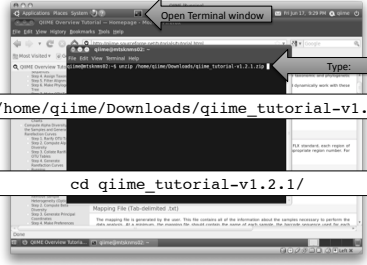
2

## Save files in the virtual machine



3

## Expand the tutorial files archive



4

## Tutorial files

Fasting_Example.fna	Fasta sequences
Fasting_Example.qual	Quality scores
Fasting_Example.sff	Binary sequence/quality file from 454 machine
Fasting_Map.txt	Meta-data for the samples in the run
custom_parameters.txt	Parameters of QIIME analysis
qiime_tutorial_commands_serial.sh	Commands to execute QIIME pipeline

5

## Quick glance at the files

```
> head Fasting_Example.fna
>FLP3FBN01ELBSX length=250 xy=1766_0111 region=1 ...
ACAGAGTCGGCTCATGCTGCCTCCCGTAGGAGTCTGGGCCGTGCTCAGT
CCCAATGTGGCCGTTTACCCTCTCAGGCCGGCTACGCATCATCGCCTTGG
TGGGCCGTTACCTCACCACACTAGCTAATGCGCCGAGGTCATCCATGTT
CACGCCTTGATGGGGCCTTAAATATACTGAGCATCGCCTCTGTATACCTA
TCCGGTTTGTAGTACCGTTTCCAGCAGTATCCCGGACACATGGGCTAGG
>FLP3FBN01EG8AX length=276 xy=1719_1463 region=1 ...
ACAGAGTCGGCTCATGCTGCCTCCCGTAGGAGTTTGGACCGTGTCTCAGT
TCCAATGTGGGGCCTTCTCTCAGAACCCCTATCCATCGAAGGCTTGGT
GGGCCGTTACCCGCCAACAACCTAATGGAACGCATCCCCATCGATGACC
GAAGTCTTAAATAGTCTTACCATCGGAA
```

6



## Sequences separated by barcode are re-labeled with new identifiers

```
>PC.634_1 FLP3FBN01ELBSX orig_bc=ACAGAGTCGGCT
new_bc=ACAGAGTCGGCT bc_diffs=0
CTGGGCCG ...
>PC.634_2 FLP3FBN01EG8AX orig_bc=ACAGAGTCGGCT
new_bc=ACAGAGTCGGCT bc_diffs=0
TTGGACCG ...
>PC.354_3 FLP3FBN01EEWKD orig_bc=AGCACGAGCCTA
new_bc=AGCACGAGCCTA bc_diffs=0
TTGGGCCG ...
>PC.481_4 FLP3FBN01DEHK3 orig_bc=ACCAGCGACTAG
new_bc=ACCAGCGACTAG bc_diffs=0
CTGGGCCGTG ...
```

13

## Define OTUs

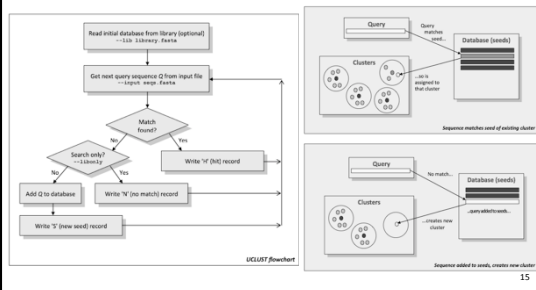
> pick\_otus\_through\_otu\_table.py ...

1. Pick OTUs with uclust at similarity of 0.97;
2. Pick a representative set with the first method;
3. Align the representative set with PyNAST
4. Assign taxonomy with RDP classifier;
5. Filter the alignment prior to tree building - remove positions which are all gaps, and specified as 0 in the lanemask
6. Build a phylogenetic tree with FastTree;
7. Build an OTU table.

14

## uclust

- <http://drive5.com/usearch/usearch3.0.html>

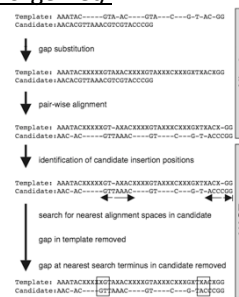


15

## PyNAST (Nearest Alignment Space Termination)

- <http://pynast.sourceforge.net/>

PyNAST: a flexible tool for aligning sequences to a template alignment.  
 Caporaso, JG.; Bittinger, K; Bushman, FD.; DeSantis, TZ.; Andersen, GL.; Knight, R.  
 January 15, 2010, DOI 10.1093/bioinformatics/btp636. Bioinformatics 26: 266-267.



16

## RDP classifier

- <http://rdp.cme.msu.edu/classifier/>

**Classes:** Well-defined fully taxonomies of bacteria

**Features:** frequencies of all  $4^8 = 65,536$  possible 8-base

**Assignment:** According to highest posterior probability

**Bootstrap significance:** Randomize input sequence 8-mers

Under conditional independence assumption the classification model is:

$$\begin{aligned}
 P(C | F_1, \dots, F_n) &= \frac{P(C)P(F_1, \dots, F_n | C)}{P(F_1, \dots, F_n)} \propto P(C)P(F_1, \dots, F_n | C) \\
 &= P(C)P(F_1 | C) \times P(F_2 | C, F_1) \times \dots \times P(F_n | C, F_1, \dots, F_{n-1}) \\
 &= P(C)P(F_1 | C) \times P(F_2 | C) \times \dots \times P(F_n | C) = P(C) \prod_{i=1}^n P(F_i | C)
 \end{aligned}$$

Wang, Q., G. M. Garrity, J. M. Tiedje, and J. R. Cole. 2007. Naive Bayesian Classifier for Rapid Assignment of rRNA Sequences into the New Bacterial Taxonomy. Appl Environ Microbiol. 73(16):5261-7.

17

## FastTree

- <http://www.microbesonline.org/fasttree/>

Price, M.N., Dehal, P.S., and Arkin, A.P. (2010) FastTree 2 -- Approximately Maximum-Likelihood Trees for Large Alignments. PLoS ONE, 5(3):e9490. doi:10.1371/journal.pone.0009490.

Heuristic phylogenetic tree reconstruction method consisting of the following steps:

1. Heuristic neighbor-joining
2. Reducing the length of the tree:
  - Nearest-neighbor interchanges
  - Subtree-prune-regraft moves
  - Distance model moves
3. Maximizing the tree's likelihood with NNIs
4. Local support values

18

## FastTree

Nearest-neighbor interchange:

Subtree-prune-regraft:

Distance model moves:  
Update the branch lengths using  
Jukes-Cantor substitution distance:  
 $-0.75 \log(1 - \frac{4}{3} d)$ ,  
where  $d$  is the proportion of  
positions that differ

19

## (some) Output files

```
wf_da/
log_[DATETIME].txt - summary of parameters and commands
uclust_picked_otus/
seqs_otus.txt - listing of sequences belonging to OTUs
rep_set/
seqs_rep_set.fasta - representative sequences for each OTU
pynast_aligned_seqs/
fasttree_phylogeny/
seqs_rep_set.tre - NEWICK format phylogeny reconstructed
seqs_rep_set_aligned.fasta - aligned representative sequences
rdp_assigned_taxonomy/
otu_table/
seqs_otu_table.txt - OTU table
seqs_rep_set_tax_assignments.txt - taxonomic assignment of OTUs
```

20

## Descriptive analyses

- Heatmaps
- Co-occurrence network
- OTU abundance charts

21

## Descriptive analysis output

- wf\_da/uclust\_picked\_otus/  
rep\_set/rdp\_assigned\_taxonomy/  
otu\_table/  
OTU\_Heatmap  
OTU\_Network  
Taxa\_Charts

22

## Diversity analyses

- alpha\_rarefaction.py ...

Rarefaction:  
- Normalization to equal number of sequences – shows the trends and if the depth is adequate  
(some) Diversity metrics:  
chao1:  $S_{obs} + n_1(n_1 - 1) / (2n_1 - 2)$   
dominance (probability of randomly drawing two individuals of the same species):  $\sum (S_i^2) / (N(N-1))$   
fisher\_alpha: solution to  $S = a * \ln(1 + n/a)$ , where  $S$  is number of species  
observed\_species: total number of unique species (richness)  
shannon: (evenness)  $-\sum p_i \log(p_i)$   
simpson:  $1 - \text{dominance}$   
singles: number of species represented by a single individual  
PD\_whole\_tree: sum of branch lengths between all representatives

23

## Alpha rarefaction results

```
> firefox wf_arare/alpha_rarefaction_plots/rarefaction_plots.html
```

24

## Beta diversity

- Compares the diversity (similarity) between two samples
- e.g. Sørensen's similarity index
- Corresponding dissimilarity (distance) is Bray-Curtis distance  $D = 1 - \beta$
- Unifrac is another popular diversity metric defined as the sum of branch lengths of phylogenetic trees that are in common between two samples:

Lozupone and Knight, UniFrac: a New Phylogenetic Method for Comparing Microbial Communities. Appl Environ Microbiol. 2005 December; 71(12): 8228-8235.

25

## Beta diversity & PCA analysis

- beta\_diversity\_through\_3d\_plots.py ...
- make\_2d\_plots.py ...
- make\_distance\_histograms.py ...

26

## Beta diversity & PCA output

- wf\_bdiv\_even146/
  - Distance\_Histograms/
  - PCoA plots
    - unweighted\_unifrac\_3d\_continuous/
    - weighted\_unifrac\_3d\_continuous/
    - unweighted\_unifrac\_3d\_discrete/
    - weighted\_unifrac\_3d\_discrete/
  - eigenvectors
    - weighted\_unifrac\_pc.txt
    - unweighted\_unifrac\_pc.txt
  - Unifrac distances
    - weighted\_unifrac\_seqs\_otu\_table\_even146.txt
    - unweighted\_unifrac\_seqs\_otu\_table\_even146.txt

27