

# otupipe

<http://drive5.com/otupipe>  
[robert@drive5.com](mailto:robert@drive5.com)

Version 1.1  
August 19, 2011

## Introduction

The otupipe script creates OTUs from next-generation sequence reads for single-region experiments such as 16S and ITS. Otupipe is a bash script, so you must have a bash shell installed. If you're using Windows, you can run otupipe under [Cygwin](#).

## Designed for 454 fixed-region sequencing

Version 1.1 is designed for reads generated by Roche 454 sequencing of a region of a gene extracted from a fixed pair of primers. Typically the gene is 16S, but this script should also work well for other regions such as ITS.

## Using otupipe with Illumina

Experiments using Illumina sequencing on the 16S gene appear to show that error rates are much higher than the rates found in genome resequencing. The reasons for this are not understood, and to the best of my knowledge no effective denoising / error-correction solution for Illumina has been published for this type of experiment. Otupipe probably works at least as well as any other solution for generating OTUs from Illumina reads. I suggest setting the MINSIZE option to a much larger value, perhaps something like 100 to 1000, and keeping only the largest OTUs, perhaps with a minimum size of 5000 to 10000.

## USEARCH

Otupipe requires [USEARCH](#), so you should install this first if you don't already have it.

## Installation

Otupipe is distributed as a tarball named something like otupipe1.1.123.tz, where 1.1.123 is the version number. To extract the files, use:

```
tar -zxvf otupipe1.1.123.tz
```

This creates a sub-directory otupipe1.0.123 with the following files:

```
otupipe.bash  
fasta_number.py
```

Before running `otupipe.bash`, you should set up your environment as explained in the following table.

Action	Description
<b>Install <a href="#">USEARCH</a>.</b>	The <code>otupipe</code> script requires <a href="#">USEARCH</a> . To install, download the binary file and install in a directory that will be accessible from your bash <code>\$PATH</code> when you run <code>otupipe</code> .
<b><code>fasta_number.py</code> in <code>\$PATH</code>.</b>	The <code>fasta_number.py</code> script must be accessible through your bash <code>\$PATH</code> when you run <code>otupipe.bash</code> . You should copy <code>fasta_number.py</code> to a directory that is in your <code>\$PATH</code> variable, or update your <code>\$PATH</code> to include the <code>otupipe</code> directory.
<b>Set <code>\$u</code>.</b>	The <code>\$u</code> environment variable must be set to the name or path of your <a href="#">USEARCH</a> binary. To verify that it is set correctly, try this command: <pre>\$u --version</pre>
<b>Set <code>\$UCHIME_REFDB</code>.</b>	The <code>\$UCHIME_REFDB</code> database must be set to the path name of the reference database to use for <a href="#">UCHIME</a> in reference database mode. You don't need to install <code>UCHIME</code> separately; <code>USEARCH</code> has its own implementation of the <code>UCHIME</code> algorithm.  For a reference database, I generally recommend the 'gold' set distributed with <a href="#">ChimeraSlayer</a> . For convenience, a copy of the gold set is available for download at <a href="http://drive5.com/otupipe">http://drive5.com/otupipe</a> , but I recommend checking the <a href="#">ChimeraSlayer page</a> to see if a more recent version is available.  Example: <pre>export UCHIME_REFDB=/public/otupipe/gold.fa</pre>

### Basic usage

The `otupipe` script requires two arguments: an input file in FASTA format containing quality-filtered reads and the name of a directory to contain the output files; this directory will be created if it does not already exist. For example:

```
otupipe.bash reads.fasta outdir [existing_otus.fasta]
```

## Output files

Three output files are created, as described in the following table.

Filename	Description
<b>otus.fa</b>	FASTA file containing one representative sequence for each OTU. Sequences are named OTU_1, OTU_2... up to the number of OTUs.
<b>chimeras.fa</b>	FASTA file containing representative sequences for chimeras found in the reads. This file is clustered at 97% by default to reduce redundancy in a similar way to OTUs. Sequences are named CHIMERA_1, CHIMERA_2 etc.
<b>readmap.uc</b>	File in UCLUST format classifying reads as OTU, chimera or unknown.

## The readmap.uc file

The readmap.uc file is generated by using the input file (reads) as a query against a database containing sequences from the otus.fa and chimeras.fa files. The UCLUST format is described in the USEARCH manual.

A UCLUST file is a tab-separated text file with one record per line. The readmap.uc file has one line for each read; either a hit (H) or no match (N) record. Most reads will match either an OTU or a chimera (H records); a small fraction of the reads will fail to match (N records). Reads that do not match are outliers that cannot be reliably classified. They might be rare biological sequences, or more likely they have errors due to PCR amplification or sequencing error.

## Example readmap.uc records

UCLUST files have ten fields. Here are some example records.

Type	Cluster	Length	PctId	Strand	Lo	Hi	Alignment	Label	Target
1	2	3	4	5	6	7	8	9	10
H	97	224	97.3	+	0	0	224M5I	FV9NWLF01B7XBX	CHIMERA_21
H	0	254	100.0	+	0	0	254M	FV9NWLF01ENNCH	OTU_1
H	180	251	97.6	+	0	0	18MI233M	FV9NWLF01EU5SK	CHIMERA_104
H	12	264	99.6	+	0	0	264MI	FV9NWLF01DYX6E	OTU_13
H	169	250	100.0	+	0	0	250M	FV9NWLF01BUPFU	CHIMERA_93
H	1	242	100.0	+	0	0	242M2I	FV9NWLF01C3TV5	OTU_2
N	*	245	*	*	*	*	*	FV9NWLF01DT3FI	*

The most relevant fields are field 1 which is H (hit) or N (no match), field 4, which gives the identity of the match ( $\geq$  the \$PCTID\_BIN parameter, which defaults to 97), field 9 which is the read label, and field 10 which is the label of the closest representative sequence (OTU\_n or CHIMERA\_n).

Typically you will need to convert the readmap.uc file to a different format for downstream tools such as [QIIME](#) or [mothur](#) that perform phylotyping and population analysis, e.g., compute alpha and beta diversities. The UCLUST format is designed to be easy to parse, so it should be straightforward to write a file format conversion script in Perl, Python etc. If you need help with this, by all means [contact me](#).

## User-settable parameters

Parameters can be changed by setting environment variables before running `otupipe.bash`. Supported parameters are described in the following table.

Parameter	Default	Description
<b>MINSIZE</b>	4	Discard clusters < MINSIZE in error correction round. Small clusters are more likely to be noise. You can reduce this value if it is important not to lose rare biological sequences, but this will typically produce more spurious OTUs.
<b>PCTID_ERR</b>	97	Identity threshold for error correction.
<b>PCTID_OTU</b>	97	Identity threshold for OTU clustering.
<b>PCTID_BIN</b>	97	Identity threshold for read classification step. Reducing this value will cause more reads to be classified (OTU or chimera).
<b>ABSKEW</b>	2	Abundance skew for <i>de novo</i> chimera filtering. See the <a href="#">UCHIME paper</a> and <a href="#">supplementary material</a> for a discussion of abundance skew.
<b>OTUTMPDIR</b>	<code>./otutmp.\$\$.\$RANDOM</code>	Name for directory to store temporary files. This directory is deleted (unless the script fails).
<b>PREFIX</b>	<code>OTU_</code>	Prefix for OTU labels.
<b>CHIMPREFIX</b>	<code>CHIMERA_</code>	Prefix for chimera labels.

For example,

```
export MINSIZE=2
export ABSKEW=4
otupipe reads.fasta out_size2_skew4
```

## Using an existing OTU file

You can optionally use an existing OTU file as follows:

```
otupipe.bash reads.fasta outdir existing_otus.fasta
```

Here, `existing_otus.fasta` is a FASTA file containing representative sequences for precomputed OTUs. For example, `existing_otus.fasta` could be centroids created by `otupipe.bash` from some other set of reads. The `existing_otus.fasta` file *must* use a different label convention to avoid collisions. For example, if you use the default labels `OTU_1`, `OTU_2`... for new OTUs, then `existing_otus.fasta` could use labels `X_1`, `X_2` and so on.

If you specify an `existing_otus.fasta` file, then:

- The `otus.fa` file will contain only new OTUs. New OTUs are created from reads that do not match the `existing_otus.fasta` file at the `PCTID_OTU` threshold, and will have `identity < PCTID_OTU` with all sequences in `existing_otus.fa`. New OTUs will be numbered 1, 2 etc.
- The `readmap.uc` file will map reads to both existing and new OTUs.