

October 13, 2014

MetAmp: a tool for Meta-Amplicon analysis User Manual

Ilya Y. Zhbannikov¹, Janet E. Williams¹, James A. Foster^{1,2,3}

³Institute for Bioinformatics and Evolutionary Studies, University of Idaho, Moscow, ID 83844-3051,

²Department of Biological Sciences, University of Idaho, Moscow, ID 83844-3051,

¹Department of Bioinformatics and Computational Biology, University of Idaho, Moscow, ID 83844-3051,

Keywords:

Bioinformatics, 16S genes, microbial studies, meta-analysis, data mining

Correspondence:

Ilya Y. Zhbannikov

Department of Bioinformatics and Computational Biology

University of Idaho

Moscow, ID 83844-30521

Email: ilyaz@uidaho.edu, foster@uidaho.edu

1 Description

MetAmp tool was developed to analyze microbial amplicon data by combining several marker regions from 16S rRNA gene. Such marker regions serve as unique identifications for species. There are nine marker regions in total in bacterial 16S rRNA gene.

MetAmp was developed by Ilya Y. Zhbannikov (ilyaz@uidaho.edu, i.zhbannikov@mail.ru), James A. Foster (foster@uidaho.edu) and Janet Williams (janetw@uidaho.edu).

2 Installation

MetAmp is easy to install. However, you must have R (www.r-project.org), Python (www.python.org) and GCC (<https://gcc.gnu.org/>).

1. Download or clone MetAmp from our GitHub repository (<http://github.com/izhbannikov/MetAmp>), for example:

```
git clone http://github.com/izhbannikov/MetAmp
```

2. In Terminal: `cd` to the MetAmp directory and execute Makefile: `make`. First, the installation script will check for required tools: R, Python and GCC. Then it builds and installs required packages and data for evaluation. You can also run `make install` to avoid checking stage.

3. Installation summary:

```
$git clone http://github.com/izhbannikov/MetAmp
```

```
$cd MetAmp/
```

```
$make
```

3 Quick start

1. In `main.R`:

- Set the path to the program directory:

```
dir_path <- " /Projects/metamp/"
```

- Set the analysis directory:

```
analysis_dir <- "analysis"
```

”Analysis” directory keeps all output files, including temporary files (folder `tmp`).

- Provide your sequence library (amplicon sequence data), for example:

```
libs <- c("data/staggered/SRR072223_V13V31_1_relabeled.fastq",  
"data/staggered/SRR072223_V35V53_1_relabeled.fastq",  
"data/staggered/SRR072223_V69V96_1_relabeled.fastq")
```

- Provide reference sequences, for example:

```
# Reference 16S gene sequences  
ref16S <- "data/16S_gold_hmc.fasta"  
refs <- c("data/16S_gold_hmc_V13V31.fasta", # V1-3  
"data/16S_gold_hmc_V35V53.fasta", # V3-5  
"data/16S_gold_hmc_V69V96.fasta") # V6-9
```

2. Then run the script `main.R` from R-environment:

```
source("main.R")
```

3.1 Test data

The following test data sets were used:

- data/even/
- data/staggered/

The even and staggered are the Human mock community pyrosequence data (SRX021555), even and staggered community, 21 microbial species and one eucariotic species.

3.2 Results

3.2.1 Clustering results

By default it is `clusters.clstr` but the name can be changed (by changing `clust_filename` parameter in `config.R` file). This file has the uc format and described at <http://drive5.com/usearch/manual/ucout.html>. Below we provide an example and copy the format description from www.drive5.com:

```
H      0      250      99.6    +      0      0      250M    read11;barcodeLabel=TCAG;    OTU_1
H      1      250      100.0   +      0      0      250M    read6;barcodeLabel=TCAG;    OTU_2
H      1      250      100.0   +      0      0      250M    read12;barcodeLabel=TCAG;   OTU_2
H      5      250      100.0   +      0      0      250M    read7;barcodeLabel=TCAG;    OTU_4
H     22      250      98.8     +      0      0      250M    read22;barcodeLabel=TCAG;   OTU_20
```

Description of uc format is directly taked from <http://drive5.com/usearch/manual/ucout.html>

Field	Description
1	Record type S, H, C or N (see table below).
2	Cluster number (0-based).
3	Sequence length (S, N and H) or cluster size (C).
4	For H records, percent identity with target.
5	For H records, the strand: + or - for nucleotides, . for proteins.
6	Not used
7	Not used
8	Compressed alignment or the symbol '=' (equals sign).
9	Label of query sequence (always present).
10	Label of target sequence (H records only).

Record	Description
H	Hit. Represents a query-target alignment
S	Centroid (clustering only)
C	Cluster record (clustering only)
N	No hit

3.2.2 OTU table

An OTU table, where rows are OTUs and columns are barcodes (by default - `otu_table.txt`).

Example:

OTUId	TCAG
OTU_8	123
OTU_9	158
OTU_2	341
...

3.2.3 Coordinates

A text file that contain coordinates for each point, including reference points (by default - `coordinates.crd`)

3.2.4 Log

A text file containing all messages during the analysis process (by default - `log.txt`)

4 MetAmp workflow

In the following sections we present key stages of our meta-amplicon analysis pipeline.

4.1 Input data

Input data can be any of the following NGS libraries: Roche 454, Ion Torrent, Illumina.

Roche 454 and Ion Torrent provide single-end sequences in special binary Sequence Flowgram Format (SFF). Before using SFF files you have to convert them into FASTQ format.

Illumina paired-end reads should be merged into longer single-end sequences.

Important! - Before analysis, you may need to perform re-labeling of your read headers: a read header should contain `barcodeLabel` (see more at www.drive5.com), for example:

```
@read1;barcodeLabel=TCAG;  
GATGAACGCTGGCGGCGTGCCTAATACATGCAAGT...AT  
+  
IIIAAAIIIIIIIIIIIIIIHHHIAAAAAAAAAA...II
```

If you use Illumina sequence data, you also have to merge overlapping paired-end reads. Later I will provide the scripts that can do these things above, but for now you can simply use scripts provided at www.drive5.com.

4.2 Data denoising

Data denoising is necessary for downstream analysis because it removes the majority of 'foreign' nucleotides, such as barcodes and primers, and low-quality regions. We suggest you to use SeqyClean tool for this operation, since it is the most comprehensive tool for sequence denoising. You can download SeqyClean here: <https://bitbucket.org/izhbannikov/seqyclean>. Using SeqyClean is simple and straightforward. Nevertheless, we provided a special script that does all the data de-noising.

4.3 Analysis workflow

1. Building a reference topology of microbial populations, where pairwise distances are computed from applying pairwise alignment of complete reference 16S gene sequences from RDP database, see Figure 1(a), top plane. We use the percent sequence identity in order to compute the distances between sequences.

2. Computing a guide "amplicon" topology for each reference 16S sequence. To do this, we extract marker regions (for example V1- V3 and V6-V9) from reference complete 16S sequences ("reference points") and place them on to the plane with the same methodology that was used for building the reference topology of complete 16S (see Figure 1(a), bottom plane, hollow green circles).
3. add empirical amplicon sequences (in fact, amplicon consensus sequences), obtained through amplicon sequencing of microbial data (Figure 1b, bottom, filled blue circles). This topology is the same empirical topology as if the empirical data would be clustered with existing methods
4. Normalized empirical topologies are formed foreach marker through mapping of each empirical point back to the reference 16S sequence topology. Mapping between the reference topology and guide "amplicon" topology is achieved with affine transformation. Guide sequences are so match with the corresponding reference 16S sequences, from where they were obtained from (see Figure 1a, top, solid green circles), and, in turn, carry the empirical amplicon sequences (or consensus amplicon sequences) back to the 16S plane.
5. Repeating stages 2-4 for each variable region.
6. Clustering and building OTUs.
7. Final statistics.

4.4 Output files

- Clustering results (`*.clstr`) - by default it is `clusters.clstr` but the name can be changed (by editing `clust_filename` parameter in `config.R` file)

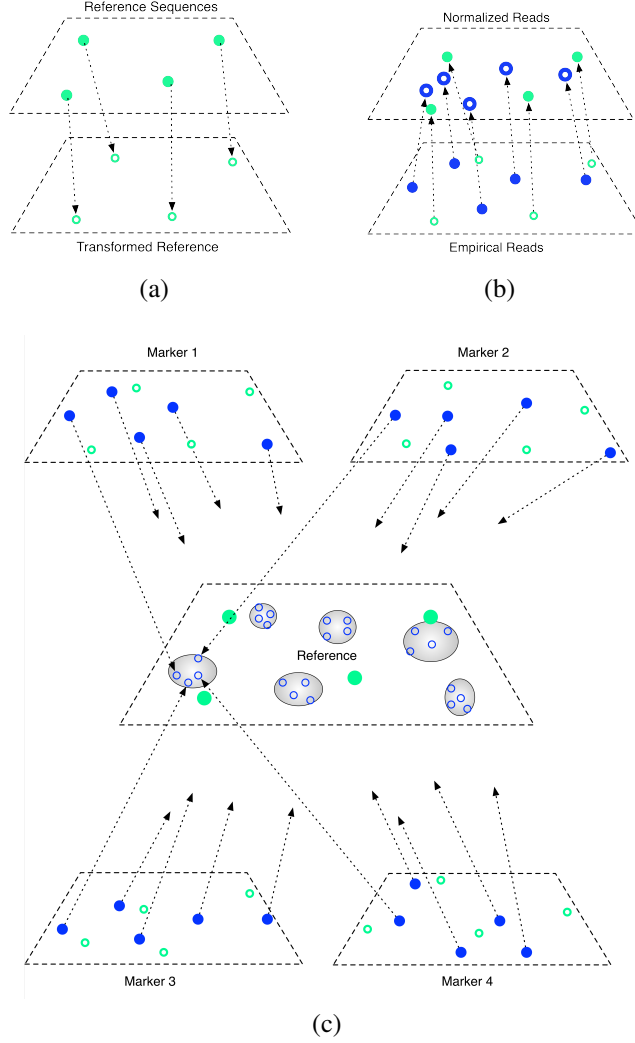


Figure 1: Illustration of “meta-amplicon” analysis algorithm. (a) Computing reference (top plane) and guide (bottom plane) topologies. Pairwise sequence dissimilarity is used to compute distances between reference sequences. The same approach is used for guide sequences, originated from corresponding whole 16S gene sequences. Reference and guide points are then placed on to the planes with multidimensional scaling. (b) The empirical amplicon sequences (hollow blue circles) are placed on to the bottom plane along with the guide circles (hollow green) circles and with the same metric. Such guide sequences are then mapped back to the reference 16S plane, carrying empirical sequences (solid blue circles). (c) Then this is repeated for each of marker (variable) region.

- An OTU table, where rows are OTUs and columns are barcodes (by default - `otu_table.txt`).
- A text file that contain coordinates for each point, including reference points (by default - `coordinates.crd`)
- A text file containing all messages during the analysis process (by default - `log.txt`)

5 Provided data

We provided some data that can be used in your projects. Notice that the amount of data is very large (even unzipped files can take several Gb). For evaluation purposes we provided small data set. Description of data uploaded to our storage is given below:

- gold folder contains 'Gold Dataset' (`gold.fa`) containing the ChimeraSlayer reference database in the Broad Microbiome Utilities (<http://microbiomeutil.sourceforge.net/>), version `microbiomeutil-r20110519`. This contains all sequences in both orientations (forward and reverse-complemented.).
- LTP folder contains data (16S sequences, extracted marker regions and distance matrices) for data extracted from All Species Living Tree Project.

LTP:

Name	Description
LTP-10271.fasta	Whole reference 16S sequences
LTP-10271_V13.fasta	reference sequences for marker 1-3
LTP-10271_V35.fasta	sequences for marker 3-5
LTP-10271_V69.fasta	sequences for marker 6-9

- Evaluation data set contains a small set of reference sequences (21 sequences) and empirical amplicon sequences (`data/even.zip` and `data/staggered.zip`). Amplicon sequences

contain sequence data from Human Mock Community. More information provided here:

<http://www.hmpdacc.org/HMMC/>

Evaluation:

Name	Description
even	Even community: SRR053818, SRR072220, SRR072239
staggered	Staggered community: SRR072221, SRR072223, SRR072237
16S_gold_hmc.fasta	contains 21 reference 16S sequences
16S_gold_hmc_V13V31.fasta	contains 21 reference marker sequences (1-3 regions)
16S_gold_hmc_V35V53.fasta	contains 21 reference marker sequences (3-5 regions)
16S_gold_hmc_V69V96.fasta	contains 21 reference marker sequences (6-9 regions)

The even and staggered are the Human mock community (HMP) pyrosequence data (SRX021555: <http://www.ncbi.nlm.nih.gov/sra?term=SRX021555>). Detailed description of these datasets (and sequencing protocols) is under the following link: <http://www.hmpdacc.org/HMMC/>

To run the program on test data open `main.R` and set the work directory (this directory should contain test data directory) and run the script `main.R` from R-environment:

```
>source("main.R")
```

Inside of the program directory we also provided some data.

5.0.1 Directory map and descriptions of folders

```
MetAmp
├── data
│   ├── gold
│   ├── LTP
│   └── RDP
├── src
│   └── Extracter
├── Packages
│   └── blastparser
└── bin
```

|_python_scripts

data - contains reference sequences, evaluation data and other data that can potentially be useful

gold - contains 'gold' reference sequences from Broad Institute. **data/LTP** - contains reference sequences and pre-computed distance matrices for > 10,000 species from All Species Living Tree Project.

src - contains source files.

Packages - contains installation files.

Packages/blastparser - a parser for blastn files.

bin - binary files (alignment and clustering programs).

python_scripts - Python scripts from www.drive5.com.

6 Parameters

You can adjust parameters in to files:

- `main.R` - common parameters, such as program directory and paths to data.
- `config.R` - advanced parameters, such as path to temporary directory, names of the output files.

Below we describe those parameters in details.

6.1 main.R

dir_path - Path to the program directory

analysis_dir - Directory that keeps all analysis data and result files

libs - data (empirical amplicon sequence libraries)

ref16S - path to the file with reference 16S sequences (file must be in FASTA format)

refs - path to the file with guide sequences, i.e. variable regions, extracted from provided file in **ref16S** (file must be in FASTA format)

6.2 config.R

usearch and **usearch7** - paths to alignment tool (we use USEARCH, see www.drive5.com)

R_LIBS - Path to installed BLASTParser package (BLASTParser parses output files from USEARCH and creates distance matrices)

clust_filename - Filename of file with final clusters

otu_table_filename - Filename of final OTU table

coord_filename - Filename of output file with NMDS coordinates of each reference and empirical point (see 'Analysis workflow')

chime_ref - Chimeric reference database

tmp_dir - A directory that contains temporary files

keep_tmp_files - Keep or not temporary files

7 Acknowledgements

This work was made possible by NIH Grants P20GM16448 and P20GM16448, and NSF Grant DBI0939454.