



Tanabe, Akifumi S.

Metabarcoding and DNA barcoding for Ecologists

Matabarcoding and DNA barcoding for Ecologists

Akifumi S. Tanabe

June 23, 2016

Table of Contents

Preface	1
Legends	3
Chapter 0	Installing softwares and preparing analysis environment 5
0.1	Installation of Claident, Assams, databases, and the other required programs 5
0.1.1	Upgrading to new version 6
0.1.2	Installing to non-default path 7
0.1.3	How to install multiple versions in a computer 8
Chapter 1	Sequencing of multiple samples by next-generation sequencers 9
1.1	PCR using tag- and adapter-jointed-primers 10
1.1.1	Decreasing costs by interim adapters 11
1.1.2	Improving sequencing quality by N 12
Chapter 2	Clustering of sequences 13
2.1	Preprocessing 13
2.1.1	In the case of Roche GS series or Life Technologies Ion PGM 13
Generating FASTQ from SFF 13	
Demultiplexing sequences 14	
Filtering low-quality sequences 16	
2.1.2	In the case of Illumina sequencers 17
Generating FASTQ from BCL 17	
Demultiplexing sequences 18	
Concatenating forward and reverse sequences 19	
Filtering low-quality sequences 19	
2.1.3	If you replicated PCR or sequencing of same template samples 20
If you replicated PCR of same template sample using same tag-jointed-primers,	
and sequenced in same sequencing run 20	
If you replicated PCR of same template sample using different tag-jointed primers,	
and sequenced in same sequencing run 21	
If you sequenced a PCR amplicon in multiple sequencing runs 21	
If you replicated PCR of same template samples using same tag-jointed primers,	
and sequenced in different sequencing runs 21	
If you replicated PCR of same template samples using different tag-jointed primers,	
and sequenced in different sequencing runs 22	

2.2	Removal of noisy and/or chimeric sequences based on sequence dereplication	23
2.3	Sequence clustering within sample	25
2.4	Sequence clustering among sample	25
2.5	Summarizing result of clustering and postprocessing	26
2.5.1	Generating OTU x Sample table	26
2.5.2	Eliminating OTUs and/or samples from table	27
2.5.3	Removal of noisy and/or chimeric sequences based on PCR replication method . . .	28
2.5.4	Removal of rarely observed OTU sequences	29
2.5.5	Dividing sequences by conserved region	29
2.5.6	Extracting ITS or SSU rRNA sequences by ITSx or Metaxa	30
2.5.7	Removal of non-target sequences	30
2.6	Submission of sequence data to DRA	31
2.6.1	Generating tab-delimited text	31
2.6.2	Generating DRA XML from tab-delimited text	32
Chapter 3	Identification of host organisms of sequences by DNA barcoding	33
3.1	Retrieving neighborhood sequences by BLAST search	33
3.2	Taxonomic assignment based on neighborhood sequences	34
3.3	Integrating multiple identification results	36
Appendix A	Installing optional programs	41
A.1	Installation of bcl2fastq	41
Appendix B	Useful terminal commands	43
B.1	Counting sequences	43
B.2	Viewing sequences	43
B.3	Compression and decompression	44
B.4	Sequence selection	44

Preface

I started the writing of this text to promote amplicon sequence assembler “Assams” and utility programs package “Claident” for sequence clustering and identification. In this text, I explain not only usage of these programs but also the methods and procedures of data collection and taxonomic identification by DNA barcoding.

Metabarcoding has been already widely used in bacterial reseaches, but it’s utility is not limited to bacteria. Metabarcoding is also applicable to soil fungi, fungi inhabiting the bodies of animals and plants, aquatic planktons, and environmental DNA emitted from macro-organisms. I explain the methods of amplification of barcoding locus from environmental DNA or metagenomes, sequencing by high-throughput sequencers (a.k.a next-generation sequencers), taxonomic assignment of nucleotide sequences (a.k.a. DNA barcoding), and observation of presence/absence of operational taxonomic units.

This text is distributed under Creative-Commons Attribution-ShareAlike 2.1 Japan License. You can copy, redistribute, display this text if you designate the authorship. You can also modify this text and distribute the modified version if you designate the authorship and apply this license or compatible license to the modified version. You can read the license at the following URL.

<http://creativecommons.org/licenses/by-sa/2.1/jp/>

You can also ask about this license to Creative-Commons, 171 Second Street, Suite 300, San Francisco, California 94105, USA.

I hope that this text helps you. I am grateful to Dr. Hirokazu Toju (Graduate Shool of Human and Environmental Sciences, Kyoto University), Dr. Satoshi Nagai (National Research Institute of Fisheries Science, Japan Fisheries Research and Education Agency), and you.

Legends

In this text, the input commands to terminals and display outputs are described as below.

```
# comments
> command argument1 \
argument2 \
argument3 ↓
output of command
> command argument1 argument2 argument3 ↓
output of command
```

In the above example, the same commands `command argument1 argument2 argument3` were executed twice. The outputs `output of command` were displayed after execution. The characters between `#` and line feed were comments and needless to input. `>` and space of line head indicate the prompt of terminal. Do not type these characters. `↓` means the end of input commands and arguments and needless to input, but you need to type Enter key to input line feed. I use line feed within commands or arguments for viewability. Such line feed is led by `\`. Therefore, the line feeds led by `\` do not mean the end of commands or arguments, or designation to input Enter key. Involuntary line feeds may be generated by word wrap function depending on your read environment, but do not mean the end of commands or arguments, or designation to input Enter key.

The file content is shown as below in this text.

```
| The content of first line
| The content of second line
```

`|` and space of line head indicate the line head in the file, do not exist in the file and needless to input these characters. This code is written to help you to distinguish true line feeds and involuntary line feeds.

Chapter 0

Installing softwares and preparing analysis environment

In this text, I assume Debian GNU/Linux jessie (hereafter Debian) or Ubuntu Linux 14.04 LTS (hereafter Ubuntu) as operating system. If you use Windows PC, please install Debian or Ubuntu. You can use CD, DVD or USB memory to boot installer. If your PC has only one storage device, you need to reduce Windows partition by using partition resizer software such as EaseUS Partition Master or using a partition resize function contained in the installer. You can also use newly added internal storage devices or external storage devices connected by USB. There are several variations of Ubuntu, and I recommend Xubuntu rather than normal Ubuntu.

Debian and Ubuntu can be installed to Mac. If there is no enough space, you need to resize OSX partition with the aid of Disk Utility or add storage device. The rEFIt or rEFInd boot selector may be required to boot Debian, Ubuntu or the installer of them on Mac. If you install rEFIt or rEFInd to your Mac, you can boot the installer of Debian or Ubuntu from CD, DVD or USB memory. Do not delete existing partition of OSX. If you have enough free space, you don't need to use Disk Utility to resize existing partition. You can install Debian or Ubuntu to external storage devices on Mac.

I assume Intel64/AMD64 (x86_64) CPU machine as analysis environment. The other CPU machine can be used for analysis, but you need to solve problems by yourself. The 64 bits version of Debian or Ubuntu is also required because 32 bits version cannot use large memory.

0.1 Installation of Claident, Assams, databases, and the other required programs

Run the following commands in terminal or console as the user that can use sudo. Then, all of the required softwares will be installed. The installer will ask password to you when sudo is used.

```
> mkdir -p ~/workingdirectory ↓
> cd ~/workingdirectory ↓
> wget https://www.claident.org/installClaident_Debian.sh ↓
> sh installClaident_Debian.sh ↓
> wget https://www.claident.org/installOptions_Debian.sh ↓
> sh installOptions_Debian.sh ↓
> wget https://www.claident.org/installDB_Debian.sh ↓
> sh installDB_Debian.sh ↓
> wget https://www.claident.org/installUCHIMEDB_Debian.sh ↓
> sh installUCHIMEDB_Debian.sh ↓
> cd .. ↓
> rm -r workingdirectory ↓
```

By default, the softwares will be installed to `/usr/local`. In the installation, you will see `Permission denied` error and the installer ask password to you. If the installer continue after password input, you don't need to care about the error. The installer try to install without `sudo` at first and the installation output the above error. Then, the installer try to install using `sudo`.

If you need proxy to connect the internet, execute the following commands to set environment variables before execution of the installer.

```
> export http_proxy=http://server.address:portnumber ↓
> export https_proxy=$http_proxy ↓
> export ftp_proxy=$http_proxy ↓
> export all_proxy=$http_proxy ↓
```

If the proxy requires username and password, execute the following commands instead of the above commands.

```
> export http_proxy=http://username:password@server.address:portnumber ↓
> export https_proxy=$http_proxy ↓
> export ftp_proxy=$http_proxy ↓
> export all_proxy=$http_proxy ↓
```

0.1.1 Upgrading to new version

If you want to upgrade all of the softwares and the databases, run the same commands as initial installation. By this procedure, Assams, Claident, PEAR, VSEARCH, Metaxa and ITSx will be installed to `/usr/local`, and NCBI BLAST+, BLAST databases for molecular identification, taxonomy databases and the other required programs will be installed to `/usr/local/share/claident`. NCBI BLAST+ and BLAST databases used by Claident can co-exist system wide installation of NCBI BLAST+ and BLAST databases.

You can disable a part of upgrade like below.

```
> mkdir -p ~/workingdirectory ↓
> cd ~/workingdirectory ↓
# disable upgrade of Assams
> touch .assams ↓
```

```

# disable upgrade of Claident
> touch .claident ↓
# disable upgrade of PEAR
> touch .pear ↓
# disable upgrade of VSEARCH
> touch .vsearch ↓
# disable upgrade of NCBI BLAST+
> touch .blast ↓
# execute upgrade > wget https://www.claident.org/installClaident_Debian.sh ↓
> sh installClaident_Debian.sh ↓
# disable upgrade of HMMer
> touch .hmmmer ↓
# disable upgrade of MAFFT
> touch .mafft ↓
# disable upgrade of Metaxa
> touch .metaxa ↓
# disable upgrade of ITSx
> touch .itsx ↓
# execute upgrade > wget https://www.claident.org/installOptions_Debian.sh ↓
> sh installOptions_Debian.sh ↓
# disable upgrade of ‘overall’ BLAST and taxonomy databases
> touch .overall ↓
# disable upgrade of ‘prokaryota’ BLAST and taxonomy databases
> touch .prokaryota ↓
# execute upgrade > wget https://www.claident.org/installDB_Debian.sh ↓
> sh installDB_Debian.sh ↓
# disable upgrade of ‘rdp’ reference database for chimera detection
> touch .rdp ↓
# disable upgrade of ‘unite’ reference databases for chimera detection
> touch .unite ↓
# execute upgrade > wget https://www.claident.org/installUCHIMEDB_Debian.sh ↓
> sh installUCHIMEDB_Debian.sh ↓
> cd .. ↓
> rm -r workingdirectory ↓

```

0.1.2 Installing to non-default path

If you install the softwares based on the above procedure, the softwares will be installed to `/usr/local`. The executable commands will be installed to `/usr/local/bin`. You can change these install path for coexistence with the other programs such as older versions like below.

```

> mkdir -p ~/workingdirectory ↓
> cd ~/workingdirectory ↓
> export PREFIX=install_path ↓
> wget https://www.claident.org/installClaident_Debian.sh ↓
> sh installClaident_Debian.sh ↓
> wget https://www.claident.org/installOptions_Debian.sh ↓
> sh installOptions_Debian.sh ↓
> wget https://www.claident.org/installDB_Debian.sh ↓
> sh installDB_Debian.sh ↓
> wget https://www.claident.org/installUCHIMEDB_Debian.sh ↓
> sh installUCHIMEDB_Debian.sh ↓
> cd .. ↓

```

```
> rm -r workingdirectory ↓
```

In this case, the following commands need to be executed before analysis.

```
> export PATH=install_path/bin:$PATH ↓
```

You can omit above command if the above command is added to `~/.bash_profile` or `~/.bashrc`.

0.1.3 How to install multiple versions in a computer

If you install Claident and the other softwares to default install path of a computer to which Claident was already installed, all softwares will be overwritten. As noted above, multiple versions of Claident can coexist if you install Claident to non-default path. Note that a configuration file `.claident` placed at a home directory of login user (`/home/username`) or `/etc/claident` cannot coexist at the same path. You need to replace this file before changing the version of Claident. The configuration file at the home directory of login user will be used preferentially. To use multiple version, I recommend to make user account for each version and to install Claident to the home directory of each user. Then, the version of Claident can be switched by switching login user.

Chapter 1

Sequencing of multiple samples by next-generation sequencers

In this chapter, I explain brief overview of tagged multiplex sequencing method by Roche GS series sequencers, Ion PGM and Illumina MiSeq. These sequencers can read over-400bp contiguously and are suitable for metabarcoding and DNA barcoding. Note that MiSeq requires concatenation of paired-end reads. Therefore, PCR amplicons should be 500bp or shorter (400bp is recommended) in order to concatenate paired-end reads. Forward and reverse reads can be analyzed separately, but I cannot recommend such analysis because reverse reads are usually low quality.

The next-generation sequencers output extremely large amount of nucleotide sequences in single run. Running costs of single run is much higher than Sanger method-based sequencers. To use such sequencers efficiently, multiplex sequencing method was developed. Multiplex identifier tag sequences are added to target sequences to identify the sample of origin, and the multiple tagged samples are mixed and sequenced in single run in this method. This method can extremely reduce per-sample sequencing costs. Multiplex identifier tag is also called as “barcode”. However, nucleotide sequence for DNA barcoding is called as “barcode sequence”. This is very confusing and “multiplex identifier tag” is too long. Thus, I call multiplex identifier tag sequence as just “tag” in this text. Please notice that tag is often called as “index”.

In the following analysis, chimera sequences constructed in PCR and erroneous sequences potentially causes misinterpretation of analysis results. If multiple PCR replicates are prepared, tagged and sequenced separately, shared sequences among all replicates can be considered as nonchimeric and less erroneous. This is because there are huge number of sequence combinations and joint points but no error sequence pattern is only one for one true sequence and nonchimeric and no error sequences likely to be observed at all replicates. Program cannot remove chimeras and errors enough but we can expect that the combination of PCR replicates and program improves removal efficiency of chimeras and errors. After removal of chimeras and errors, the number of sequences of PCR replicates can be summed up and used in subsequent analysis.

1.1 PCR using tag- and adapter-jointed-primers

In order to add tag to amplicon, PCR using tag-jointed primer is the easiest way. This method requires a set of tag-jointed primers. In addition, library preparation kits for next-generation sequencers usually presume that the adapter sequences specified by manufacturers are added to the both end of target sequences. Thus, the following tag- and adapter-jointed primer is used for PCR.

■ 5' — [adapter] — [tag] — [specific primer] — 3'

If this kind of primers are used for the both forward and reverse primers, the following amplicon sequences will be constructed.

■ 5' — [adapter-F] — [tag-F] — [specific primer-F] — [target sequence] — [specific primer-R (reverse complement)] — [tag-R (reverse complement)] — [adapter-R (reverse complement)] — 3'

In the case of single-end read, tag-F leads specific primer-F and target sequence in the sequence data.

The supplement of Hamady *et al.* (2008) may be useful for picking tag sequences. In the case of single-end sequencing, 3'-side tag is not required, and tagless primer can be used for PCR. In the case of paired-end sequencing, single index (tag) can be applied, but dual index (tag) is recommended for detecting unlikely tag combinations which means that forward and reverse sequences are mispaired.

Using above primer sets for PCR, primers anneal to templates in "Y"-formation, and the amplicon sequences which have tags and adapters for both ends will be constructed. Then, the amplicon solutions are mixed in the same concentration and sequenced based on manufacturer's protocol. Spectrophotometer (including Nanodrop) is inappropriate for the measurement of the concentration of solution because measurement of dsDNA using spectrophotometer is likely to be affected by the other contaminants. I recommend Qubit (ThermoFisher) for measurement of dsDNA concentration. Quantitative PCR-based method can also be recommended but it's expensive and more time-consuming.

Primer annealing position sequence can also be used for recognizing the sample of origin. Therefore, the sequences of multiple loci, for example plant *rbcL* and *matK*, from same sample set tagged by same tag set can be multiplexed and sequenced. Of course, the sequences of multiple loci can also be recognized by themselves. Smaller number of cycles and longer extension time were recommended for PCR. Because the required amount of DNA for sequence sample preparation is not so high, the larger number of cycles of PCR amplification is not needed. The larger number of cycles and shorter extension time generates more incompletely extended amplicon sequences and the incompletely extended amplicon sequences are re-extend using different template sequences in next cycle. Such sequences are called as "chimeric DNA". Chimeric DNAs causes a discovery of non-existent novel species or a overestimation of species diversity. To reduce chimeric DNA construction, using high-fidelity DNA polymerase such as Phusion (Finnzymes) or KOD (TOYOBO) is effective. Stevens *et al.* (2013) reported that slowing cooling-down from denaturation temperature to annealing temperature reduced chimeric DNA construction. If your thermal cycler can change cooling speed, slowing cooling-down from denaturation temperature to annealing temperature can be recommended. Chimeric DNA sequences can also be eliminated by computer programs after sequencing. Because chimera removal by programs is incomplete and the nonchimeric

sequences shrink, we cannot do better than reduce chimeric DNA construction.

In the case of hardly amplifiable templates, using Ampdirect Plus (Shimadzu) for PCR buffer or crushing by homogenizer or beads before DNA extraction is recommended. Deep freezing before crushing can also be recommended. Removal of polyphenols or polysaccharides might be required if your sample contain those chemicals. If PCR amplification using tag- and adapter-jointed-primers fail, try two-step PCR that consist of primary PCR (20–30 cycles) using primers without tags and adapters, purification of amplicons by ExoSAP-IT, and secondary PCR (5–10cycles) using amplicons of primary PCR as templates and tag- and adapter-jointed-primers.

1.1.1 Decreasing costs by interim adapters

Tag- and adapter-jointed-primers are very long and expensive. In addition, we need to buy tag- and adapter-jointed-primers for each locus. To reduce cost of tag- and adapter-jointed-primers, interim adapter-jointed primers and two-step PCR is useful. The following primer set is used in primary PCR.

■ 5' — [interim adapter] — [specific primer] — 3'

This PCR product have interim adapter sequences at the both ends. This PCR product is used as template in secondary PCR after purification. The following primer set is used in secondary PCR.

■ 5' — [adapter specified by manufacturer] — [tag] — [interim adapter] — 3'

This two-step PCR enables us to reuse secondary PCR primers. However, this two-step PCR may increase PCR errors and PCR amplification biases, and decrease target sequence lengths. Note that final PCR product is constructed as the following style.

■ 5' — [adapter-F specified by manufacturer] — [tag-F] — [interim adapter-F] — [specific primer-F] — [target sequence] — [specific primer-R (reverse complement)] — [interim adapter-R (reverse complement)] — [tag-R (reverse complement)] — [adapter-R (reverse complement) specified by manufacturer] — 3'

Illumina の Nextera XT Index Kit を用いた multiplex 法 (Illumina corporation, 2013) はこれをキットで実現しているものです。この方法で 2 つのタグ配列と両側からの解読を行う場合、最初に中間アダプター配列 F の後ろから、つまり本来のプライマー配列 F から読み始め、解読したい配列方向へ解読します。次に、中間アダプター配列 R 逆相補の後ろからタグ配列 R 逆相補を解読します。3 番目にアダプター配列 F の後ろからタグ配列 F を解読し、方向を反転させて、最後に中間アダプター配列 R の後ろから本来のプライマー配列 R、および解読したい配列 (逆ストランド) を解読します。それぞれが R1、R2、R3、R4 を含むファイル名で保存されます。R1、R2、R3 は同じ strand で、R4 だけ逆 strand になります。R2 はタグ配列 R 逆相補で、R3 がタグ配列 F です。ややこしいので注意が必要です。なお、上記のキットでは R1 と R4 のシーケンスプライマーは中間アダプター配列をターゲットとするものになっているため、本来のプライマー配列が R1 と R4 のデータ配列に入ってしまう、解読したい領域が 500bp 程度ある場合には長さが不足する可能性があります。本来のプライマーをシーケンスプライマーとして利用することでデータ配列にプライマー部分が入らないようすることも可能です。ただし、この方法では後述する N の挿入によるシーケンス品質の改善ができません。

1.1.2 Improving sequencing quality by N

Illumina では、フローセル上の蛍光を光学センサーで読み取っていますが、メタゲノムを鋳型としたユニバーサルプライマーでの PCR 増幅産物は配列が似ているため、フローセル上で DNA が近接していると区別が困難になります。また、読み始めに「ほとんどの配列が A であるため発光点が少なく真っ暗になる」と、異常と判断して解読を停止してしまうことがあるとのことです。ゲノムショットガンや RNA-seq では近接している DNA が非常によく似ている可能性がずっと低いために区別できる上、真っ暗になることもないのですが、同じ遺伝子座の増幅産物 (しかも読み始めがほとんど変異のないプライマー領域) だとまずいわけです。そこで、本来のプライマー配列と中間アダプター配列の間にわざと NNNNNN を入れてしまうことにします。そうすると、この部分の解読時に近接した DNA も区別することができ、真っ暗になることもないため、シーケンス品質が改善するとのことです (Nelson *et al.*, 2014)。MiSeq や HiSeq をお使いの方はお試し下さい。さらに、この長さをサンプルごとに変えてしまうと、解読するゲノム上の位置がずれるため、全く同じ配列が隣接していても区別できるようになるためにシーケンス品質が改善するそうです (Fadrosh *et al.*, 2014)。これらの対処法により、PhiX 標準配列を混ぜる量を減らすことができるので、一度に得られるデータ量も増加するはずです。

Chapter 2

Clustering of sequences

2.1 Preprocessing

GS シリーズシーケンサーや Ion PGM では、最終出力として*.sff というファイルが得られます。Illumina MiSeq の場合は*.fastq となります。ここでは、これらのファイルに含まれる配列をサンプルごとに分別 (demultiplex) して別々のファイルに分けた上で、低品質の配列の除去を行う方法を説明します。どの機種でもメーカーのプログラムの機能による demultiplex は品質の値を見ていなかったりするため、Claident のコマンド `clsplitseq` を利用して demultiplex を行います。以下ではコマンドは全てターミナルかコンソールで実行して下さい。基本的なターミナルの使い方の知識は持っているという前提で話を進めます。ターミナルの操作に慣れていない方は、最低でも付録 B の内容を理解できるようになっておいて下さい。

2.1.1 In the case of Roche GS series or Life Technologies Ion PGM

Generating FASTQ from SFF

まずは SFF ファイルから FASTQ 形式のファイルを作成します。

```
$ sff_extract -c SFF ファイル名
```

-c オプションを付けることで、読み始めの TCAG が除去されてファイルが作成されます。もし後述するタグ配列の冒頭に TCAG を付けている場合には、このオプションは使用しないで下さい。SFF ファイルが `HOGHEHOG.sff` だとすると、`HOGHEHOG.fastq` という FASTQ 形式のファイルとして塩基配列は保存されます。`HOGHEHOG.xml` というファイルもできますが、ここでは使いませんので削除して構いません。出力される塩基配列はタグ配列で始まり、解読開始側プライマー配列、解読対象の塩基配列データと続き、反対側プライマー配列 (の逆相補配列) で終わります。ただし、全域が読めるとは限りませんので、途中で読み終わっている配列も含んでいます。なお、ここでは `sff_extract` コマンドを用いていますが、読み始めの TCAG を除いた塩基配列の FASTQ 形式ファイルが得られれば、シーケンサーメーカーから提供される純正のコマンドを使っても構いません。逆にタグ配列の冒頭に TCAG を追加していただいても構いませんが、TCAG 部分の品質値まで厳格にチェックされます。

Demultiplexing sequences

次に、タグ配列とプライマー配列に基づいてサンプルごとに配列を別々のファイルに分離 (demultiplex) します。事前に以下のようなタグ配列を記した FASTA ファイル、解読開始側プライマー配列を記した FASTA ファイルがそれぞれ必要です。

```
>タグ ID
タグ配列
>examplesample1
ACGTACGT

>プライマー ID
プライマー配列
>exampleprimer1
ACGTACGTACGTACGTACGT
```

タグ配列には縮重コード表記は使えませんが、プライマー配列には縮重コード表記を利用可能です。タグ配列、プライマー配列ともに複数の配列を記述しておくことも可能です。なお、中間アダプターを用いたりした場合は、プライマー配列ファイルは以下のような内容で作成して下さい。

```
>プライマー ID
中間アダプター配列プライマー配列
```

つまり、プライマー配列ファイルには、タグ配列と解読対象の配列の間にある配列を記述して下さい。

これらのファイルが用意できたら、以下のコマンドで塩基配列をサンプルごとに別ファイルへ分割します。

```
$ clsplitseq \
--runname=ラン ID \
--tagfile=タグ配列ファイル \
--primerfile=プライマー配列ファイル \
--minqualtag=27 \
--numthreads=使用する CPU 数 \
入力ファイル \
出力フォルダ
```

ラン ID はランごとに異なるように適当に決めて指定して下さい。通常はシーケンサーのメーカーが適当に付けているはずなので、それを記述します。ファイル名やファイル中の配列名などに含まれていますが、メーカーによってルールが異なりプログラムから自動的に取得することが面倒なため、ユーザー側で指定する必要があります。--minqualtag はタグ配列部分の品質値の下限を指定するオプションです。この値より品質値が低い塩基がタグ配列に含まれる場合は、その配列は出力されません。品質値 27 という値は、Roche GS シリーズシーケンサーの読み取りミスをクラスタリング・コンセンサス配列作成によって除去するのに適した値です (Kunin *et al.*, 2010)。他機種では別の値の方が良い可能性があります。品質値 30 がかなり広く使用されていますので、30 にしてもいいでしょう。

出力される配列からは、タグとプライマーの配列部分は除去されています。タグ配列部分の認識は厳密一致です。1 塩基の読み間違いを許すなどのオプションはありません。プライマー配列の認識は、Needleman-Wunsch アルゴリズムを用いたアライメントを行い、14% まで不一致を許容します (閾値は変更可能)。出力フォルダには、ラン ID__タグ ID__プライマー ID.fastq.gz という名前の塩基配列ファイルが保存されます。clsplitseq では並列処理によって多数の CPU を活かした高速化が可能です。CPU が 4 個の場合、--numthreads オプションに 4 を指定すれば最も高速に処理ができるはずです。ただし、OS やディスクへの書き込み速度がボトルネックになって思うように速度が出ない

こともあります。なお、出力ファイルは GZIP で圧縮されていますので、GZIP 圧縮 FASTQ に対応していないプログラムで読み込むには `gzip -d` などでも圧縮を解除しておく必要があります。これ以降の Claident のコマンドは対応しているためそのまま使えます。

DDBJ Sequence Read Archive (DRA) など、次世代シーケンサーの生データを登録するサイトには、ここで生成された GZIP 圧縮 FASTQ を登録します。

■数ラン分に分けて大量サンプルをシーケンスした場合 `clsplitseq` による処理を複数回行います。ただし、`clsplitseq` は出力先フォルダが既に存在しているとエラーになり処理できません。そこで、2 回目以降は `--append` オプションを付けて実行します。例えば以下のようにして下さい。

```
$ clsplitseq \  
--runname=ラン ID \  
--tagfile=タグ配列ファイル \  
--primerfile=プライマー配列ファイル \  
--minqualtag=27 \  
--numthreads=使用する CPU 数 \  
入力ファイル 1 \  
出力フォルダ  
$ clsplitseq \  
--runname=ラン ID \  
--tagfile=タグ配列ファイル \  
--primerfile=プライマー配列ファイル \  
--minqualtag=27 \  
--numthreads=使用する CPU 数 \  
--append \  
入力ファイル 2 \  
出力フォルダ
```

■異なる長さのタグが混在している場合について `clsplitseq` では、タグの長さは全て同一であることを仮定しています。その仮定を使うことで処理を高速化しています。そのため、タグの長さが不統一の場合は、長さの異なるタグごとに別々の FASTA ファイルを準備し、以下のようにして複数回 `clsplitseq` による処理を行って下さい。

```
$ clsplitseq \  
--runname=ラン ID \  
--tagfile=タグ配列ファイル 1 \  
--primerfile=プライマー配列ファイル \  
--minqualtag=27 \  
--numthreads=使用する CPU 数 \  
入力ファイル \  
出力フォルダ  
$ clsplitseq \  
--runname=ラン ID \  
--tagfile=タグ配列ファイル 2 \  
--primerfile=プライマー配列ファイル \  
--minqualtag=27 \  
--numthreads=使用する CPU 数 \  
--append \  
入力ファイル \  
出力フォルダ
```

■リバースプライマー配列の認識と除去 これまでの手順では、リバースプライマー配列がアニールする部分の配列は除去されていません。これもできれば除去した方が良いでしょう。それには、まず以下のような内容のリバースプライ

マー配列のファイルを用意します。

```
>プライマー ID
プライマー配列
>exampleprimer1
TCAGTCAGTCAGTCAGTCAG
```

配列は複数記述できますが、フォワードプライマーと 1 対 1 対応していることを仮定していますので注意して下さい。フォワードプライマーと配列数が異なるとエラーになります。配列順序も対応していることを仮定していますので、リバースプライマーは同一でもフォワードプライマーが異なるサンプルがある場合、それぞれ異なるプライマー ID を割り当てて両方記述する必要があります。

ファイルが用意できたら、`clsplitseq` を以下のように実行して下さい。

```
$ clsplitseq \
--runname=ラン ID \
--tagfile=タグ配列ファイル \
--primerfile=プライマー配列ファイル \
--reverseprimerfile=リバースプライマー配列ファイル \
--reversecomplement \
--minqualtag=27 \
--numthreads=使用する CPU 数 \
入力ファイル \
出力フォルダ
```

このコマンドでは前述のコマンドの処理に加えて、さらにリバースプライマー逆相補配列を Needleman-Wunsch アルゴリズムを用いたアライメントによって 15% まで不一致を許して探します。該当する配列が見つかったら、該当箇所の先頭以降を元の配列から除去します。該当する配列が見つからなかった場合は、他の条件を満たしていればそのまま出力ファイルに保存されます。ただし、読み間違いなどによってリバースプライマーの付く部位以降が除去しきれないこともあります。なお、`--needreverseprimer` というオプションを追加することで、リバースプライマー逆相補配列に一致する部位が見つからなかった配列は出力されなくなります。最後まで解読しきれた配列だけにしたい場合はこのオプションを付けて実行して下さい。

Filtering low-quality sequences

FASTQ 形式の塩基配列には、解読時のデータから推定された塩基の品質値がありますので、その情報に基づいて質の低い配列をこの時点である程度除去して行うことができます。以下のように `clfilterseq` を実行して下さい。

```
$ clfilterseq \
--minqual=27 \
--minquallen=3 \
--minlen=350 \
--maxlen=400 \
--maxprowqual=0.1 \
--numthreads=使用する CPU 数 \
入力ファイル \
出力ファイル
```

ここで、`--minqual` は品質値の下限値で、終末端側からこの値以上の品質値が `--minquallen` の値以上に連続して出現するまで配列を削っていきます。その上で、`--minlen` の長さ未満の配列は除去し、`--maxlen` より長い場合はこの

値になるまで終末端を削ります。残った配列において、`--minqual` 未満の配列が `--maxplowqual` の割合より多い場合はその配列は除去します。通常は出力はファイルになりますが、`--output=folder` というオプションを与えた場合はフォルダを作成して入力ファイルと同じ名前のファイルに出力します。既存のフォルダに出力する場合は `--append` を追加して下さい。

`clsplitseq` の出力フォルダ内の全配列ファイルに対して同じ処理を行う場合は、以下のようなコマンドを実行して下さい。

```
$ for f in clsplitseq の出力フォルダ/*.fastq.gz
do clfilterseq \
--output=folder \
--append \
--minqual=27 \
--minqual=3 \
--minlen=350 \
--maxlen=400 \
--maxplowqual=0.1 \
--numthreads=使用する CPU 数 \
$f \
出力フォルダ
done
```

2.1.2 In the case of Illumina sequencers

Generating FASTQ from BCL

Illumina 社のシーケンサーでは、標準で `demultiplex` する機能がありますが、タグ配列部分の品質値を一切見ていないため、タグ配列の品質が低い配列がファイルに含まれています。このため、生データである BCL ファイルから、Illumina が配布している `bcl2fastq` を用いて、`demultiplex` していない FASTQ 配列を生成し、`clsplitseq` で `demultiplex` を行います。`bcl2fastq` のインストールは付録に書きましたのでそちらを参照願います。

まず、シーケンサー付属マシンから BaseCalls フォルダを含むランデータフォルダをまるごと `bcl2fastq` をインストールした PC へコピーし、以下のコマンドを実行します。なお、BaseCalls フォルダ内に `SampleSheet.csv` が存在する場合は、これをリネームか削除しておく必要があります。

```
$ cd BaseCalls フォルダの直上フォルダ
$ configureBclToFastq.pl \
--fastq-cluster-count 0 \
--use-bases-mask Y300n,Y8,Y8,Y300n \
--input-dir BaseCalls フォルダ \
--output-dir 出力フォルダ
$ make -j4
```

`--fastq-cluster-count 0` でファイルの自動分割機能を無効にし、`--use-bases-mask Y300n,Y8,Y8,Y300n` でフォワード側配列を 300bp (末端 1bp は切り捨て)、インデックス 1 (タグ配列 R 逆相補) を 8bp、インデックス 2 (タグ配列 F) を 8bp、リバース側配列を 300bp (末端 1bp は切り捨て) とし、それぞれ `*R1_001.fastq.gz`、`*R2_001.fastq.gz`、`*R3_001.fastq.gz`、`*R4_001.fastq.gz` として出力するよう設定しています。これは Dual Index (インデックスは各 8bp) で両側から各 300bp ずつ解読する (PE300) 設定でシーケンスした場合ですので、異なる

る設定でシーケンスした場合は値を適宜変更して下さい。make -j4 を実行すると、4 つの CPU を使用して BCL から FASTQ の生成が行われます。生成された FASTQ は GZIP というプログラムで圧縮されたものになります。ファイル名の末尾の拡張子 .gz はそれを示すものです。Claident は GZIP で圧縮された FASTQ をそのまま扱えるため、圧縮したままで構いません。

Demultiplexing sequences

FASTQ ファイルができたなら、タグ配列とプライマー配列に基づいてサンプルごとに配列を別々のファイルに分離 (demultiplex) します。事前に以下のようなタグ配列を記した FASTA ファイル、解読開始側プライマー配列を記した FASTA ファイルがそれぞれ必要です。ペアエンドの場合は、当然リバース側の配列を記したファイルも必要です。

```
>タグ ID
タグ配列
>examplesample1
ACGTACGT

>プライマー ID
プライマー配列
>exampleprimer1
ACGTACGTACGTACGTACGT
```

タグ配列には縮重コード表記は使えませんが、プライマー配列には縮重コード表記を利用可能です。タグ配列、プライマー配列ともに複数の配列を記述しておくことも可能ですが、フォワードプライマーとリバースプライマーは配列数も配列順序も対応していることを仮定しています。リバースプライマーは同一でもフォワードプライマーが異なるサンプルがある場合、それぞれ異なるプライマー ID を割り当てて両方記述する必要があります。タグに関しても同様です。なお、プライマーの前に NNNNNN を付加した場合は、それも含めてプライマーを記述して下さい。

これらのファイルが用意できたら、以下のコマンドで塩基配列をサンプルごとに別ファイルへ分割します。

```
$ clsplitseq \
--runname=ラン ID \
--index1file=インデックス 1 配列 (タグ配列 R 逆相補) ファイル \
--index2file=インデックス 2 配列 (タグ配列 F) ファイル \
--primerfile=プライマー配列ファイル \
--reverseprimerfile=リバースプライマー配列ファイル \
--minqualtag=30 \
--numthreads=使用する CPU 数 \
入力ファイル 1 \
入力ファイル 2 \
入力ファイル 3 \
入力ファイル 4 \
出力フォルダ
```

入力ファイルは *R1_001.fastq.gz、*R2_001.fastq.gz、*R3_001.fastq.gz、*R4_001.fastq.gz の順で記述して下さい。--index1file は Illumina が言うところのインデックス 1、つまりタグ配列 R 逆相補を指定し、--index2file には Illumina が言うところのインデックス 2、つまりタグ配列 F を指定して下さい。逆にしないように注意して下さい。処理が終わったら、各ファイルの配列数を確認し、Illumina 純正のプログラムで demultiplex したファイルの配列数と比較して下さい。正しく処理できていれば、全て Illumina 純正プログラムで処理した場合より少なくなっているはずです。

出力されるファイルはラン ID_タグ ID_プライマー ID.forward.fastq.gz とラン ID_タグ ID_プライマー

ID.reverse.fastq.gz になります。DDBJ Sequence Read Archive (DRA) など、次世代シーケンサーの生データを登録するサイトには、ここで生成された GZIP 圧縮 FASTQ を登録します。なお、DRA への登録では配列の長さが揃っているかどうかを指定する箇所がありますが、長さが可変のプライマー部位の配列を除去していますから揃っていません (プライマーセットが 1 つであってもアライメント時にギャップが入る可能性があるため揃わない) ので誤って揃っていると指定しないようにして下さい。

Concatenating forward and reverse sequences

フォワード配列とリバース配列にオーバーラップがある場合には、PEAR (Zhang *et al.*, 2014) を用いて連結します。demultiplex した配列ファイルのあるフォルダで、以下のコマンドを実行して下さい。

```
$ pear \  
-p 0.0001 \  
-u 0 \  
-j 使用する CPU 数 \  
-f ラン ID__タグ ID__プライマー ID.forward.fastq.gz \  
-r ラン ID__タグ ID__プライマー ID.reverse.fastq.gz \  
-o ラン ID__タグ ID__プライマー ID
```

正常に処理が終わると、以下のファイルができます。

ラン ID__タグ ID__プライマー ID.assembled.fastq 連結された配列
ラン ID__タグ ID__プライマー ID.unassembled.forward.fastq 連結されなかったフォワード配列
ラン ID__タグ ID__プライマー ID.unassembled.reverse.fastq 連結されなかったリバース配列
ラン ID__タグ ID__プライマー ID.discarded.fastq 検定で一致していると判定されなかった配列

この後で使用するのはラン ID__タグ ID__プライマー ID.assembled.fastq だけです。また、出力は圧縮されていないため、ディスク容量の節約のためにも圧縮しておいた方がいいかもしれません。

なお、フォルダ内の全てのファイルに対して PEAR による連結を連続して行うには、以下のようなコマンドをお使い下さい。

```
$ for f in $(ls *.forward.fastq.gz | grep -P -o '^[^\.]+'  
do pear \  
-p 0.0001 \  
-u 0 \  
-j 使用する CPU 数 \  
-f $f.forward.fastq.gz \  
-r $f.reverse.fastq.gz \  
-o $f  
done
```

Filtering low-quality sequences

FASTQ 形式の塩基配列には、解読時のデータから推定された塩基の品質値がありますので、その情報に基づいて質の低い配列をこの時点である程度除去して行うことができます。以下のように `clfilterseq` を実行して下さい。

```
$ clfilterseq \  
--minqual=30 \  
--maxplowqual=0.1 \  
--numthreads=使用する CPU 数 \  
入力ファイル \  
出力ファイル
```

このコマンドでは、入力配列において `--minqual` 未満の配列が `--maxplowqual` の割合より多い場合はその配列を除去します。通常は出力はファイルになりますが、`--output=folder` というオプションを与えた場合はフォルダを作成して入力ファイルと同じ名前のファイルに出力します。既存のフォルダに出力する場合は `--append` を追加して下さい。Illumina でペアエンドを連結した場合は、両端は品質値が高く、重複部分も一致していれば品質値が高くなりますので、末端から品質値の低い配列を削るのはあまり意味がありません。連結したペアエンドではこのように品質値が低い配列が一定量を超えた配列を除去するのがよいでしょう。もちろん、既存の FASTQ の品質確認を行うプログラムを利用されてもいいと思います。そのようなプログラムとしては、FastQC (Andrews, 2010) や PRINSEQ (Schmieder & Edwards, 2011) があります。

PEAR で連結した全配列ファイルに対して同じ処理を行う場合は、以下のようなコマンドを実行して下さい。

```
$ for f in clsplitseq の出力フォルダ/*.assembled.fastq  
do clfilterseq \  
--output=folder \  
--append \  
--minqual=30 \  
--maxplowqual=0.1 \  
--numthreads=使用する CPU 数 \  
$f \  
出力フォルダ  
done
```

片側だけのシーケンスデータや、連結していないデータの場合は、Roche GS シリーズシーケンサや Ion PGM の場合と同様の処理を行えばいいと思います。第 2.1.1 節をご参照下さい。ただし、品質値や長さの閾値は適宜変更する必要があります。

2.1.3 If you replicated PCR or sequencing of same template samples

キメラ配列は PCR によって生成されますので、同一 PCR 由来配列なのか、異なる PCR 由来配列なのかをプログラムに正しく認識させなくてはなりません。そのため、同じ鋳型を複数回 PCR または複数回シーケンスした場合には、プログラムがデータを正しく扱えるようにいくつかの操作が必要になります。なお、ここでは、Nested PCR や、増幅のための PCR とタグ配列付加のための PCR はまとめて 1 回と数えることに注意して下さい。つまり、ここで「複数回 PCR」とは、「増幅産物を鋳型にしてさらに PCR で増幅する」ことではなく、「同じ鋳型を使用した PCR」を複数回行うことです。

If you replicated PCR of same template sample using same tag-jointed-primers, and sequenced in same sequencing run

実験デザインが間違っています。キメラ検出・除去が不可能です。そのデータは使えませんので破棄して下さい。

If you replicated PCR of same template sample using different tag-jointed primers, and sequenced in same sequencing run

集計表の作成後に対応するサンプル行を編集してまとめて下さい。わざと別サンプルのままにしておいて、1 つ以上もしくは全てのサンプルに出現しない OTU はキメラと判断するというのも良い考えだと思います。ただ、キメラの生成がランダムであればこれはベストな方法でしょうが、できやすいキメラというものがもしあれば、キメラ配列が全サンプルに出現する可能性があります。この方法がどの程度正確かという情報はまだ十分ではありませんので、プログラムによるキメラ検出をスキップするのはおすすめてできません。

If you sequenced a PCR amplicon in multiple sequencing runs

1 回のランでは、たまたま得られる配列が少ないことがあります。そのため同一サンプルを複数のランでシーケンスし、全ランのデータを使いたいことがあります。そのような場合は以下のように処理して下さい。

```
$ clsplitseq \  
オプション省略 \  
--runname=HOGEHOGE \  
入力ファイル 1 \  
出力フォルダ  
$ clsplitseq \  
オプション省略 \  
--runname=HOGEHOGE \  
--append \  
入力ファイル 2 \  
出力フォルダ
```

このように`--runname` オプションでラン ID を置換し、2 つの入力ファイルのラン名を同一にしてやることで、2 つのファイルからの同一サンプルからの配列が `HOGEHOGE__タグ ID__プライマー ID.fastq.gz` という 1 つのファイルに書き出されます。

Claident ではラン ID とタグ ID とプライマー ID が同じなら、同一 PCR 産物の配列として扱います。PCR が 1 回でタグ配列ファイルとプライマー配列ファイルが同じなら、当然タグ ID とプライマー ID は同じになっているはずですので、ラン ID も同じにしてしまうことで、同一 PCR 産物の配列として認識させることができるようになります。

If you replicated PCR of same template samples using same tag-jointed primers, and sequenced in different sequencing runs

複数の配列ファイルがあるはずですので、以下のように `clsplitseq` を実行して下さい。ここではファイルが 2 つの場合を示します。

```
$ clsplitseq \  
オプション省略 \  
入力ファイル 1 \  
出力フォルダ  
$ clsplitseq \  
オプション省略 \  
入力ファイル 2 \  
出力フォルダ
```

```
--append \  
入力ファイル 2 \  
出力フォルダ
```

タグ配列ファイル、プライマー配列ファイル、出力フォルダは同じものを与えて下さい。--**runname** オプションでは異なるラン ID を指定します。上記のように実行すると、2 つのランはそのまま別々のランとして扱われることになります。

Claident ではラン ID かタグ ID かプライマー ID のどれかが異なれば、異なる PCR 由来の配列として扱います。PCR が別々でも、タグ配列ファイルとプライマー配列ファイルが同じなら、当然タグ ID とプライマー ID は同じになっているはずですので、ラン ID は異なるようにしないと、異なる PCR 由来の配列として認識させることができません。また、ここでは異なるラン由来の配列は、配列ファイルも別々になっているはずですが、これ以降の処理は 1 ランでシーケンスした場合と同様になりますが、そのまま集計表を作成すると、同じ鋳型の同じ領域の配列でも、異なる PCR 由来の配列が別の行に集計されてしまいます。しかし、同じ鋳型の同じ領域の配列なのですから、集計表ではひとまとめにしてしまいたいと思います。その場合は集計表作成を以下のコマンドで行なって下さい。異なるラン ID の部分を一つのラン ID で置き換えてしまうことにより、ラン ID とタグ ID とプライマー ID が一致するようになり、単一の行にまとめて集計が行われます。

```
$ clsumclass \  
--runname=HOGHEHOGHE \  
入力ファイル \  
出力ファイル
```

If you replicated PCR of same template samples using different tag-jointed primers, and sequenced in different sequencing runs

複数の配列ファイルがあるはずですので、以下のように **clsplitseq** を実行して下さい。ここではファイルが 2 つの場合を示します。

```
$ clsplitseq \  
オプション省略 \  
--tagfile=タグ配列ファイル 1 \  
--primerfile=プライマー配列ファイル \  
入力ファイル 1 \  
出力フォルダ  
$ clsplitseq \  
オプション省略 \  
--tagfile=タグ配列ファイル 2 \  
--primerfile=プライマー配列ファイル \  
--append \  
入力ファイル 2 \  
出力フォルダ
```

プライマー配列ファイル、出力フォルダは同じものを与えて下さい。--**runname** オプションでは異なるラン ID を指定します。また、タグ配列ファイルは個別に用意して与えて下さい。ただし、同じ鋳型由来の配列に付加したタグ ID は同じにして下さい。配列が同じタグであっても、異なる鋳型由来の配列に付加したタグの ID を同一にはいけません。例えば以下の 2 ファイルを用意したとします。

```
>sample1
ACGTACGT
>sample2
ATGCATGC

>sample1
ATGCATGC
>sample2
ACGTACGT
```

この場合、sample1 の鋳型には 1 ラン目には ACGTACGT を付加し、2 ラン目では ATGCATGC を付加したことになります。sample2 の鋳型では逆になります。もしも、ほぼ同じ領域を増幅してはいるがプライマー配列は異なる場合、タグと同様に同じ ID を付けた別々のファイルを用意して与えて下さい。

Claident ではラン ID かタグ ID かプライマー ID のどれかが異なれば、異なる PCR 由来の配列として扱います。ラン ID とタグ ID が異なれば、当然異なる PCR 由来の配列として扱われるのですが、同じ鋳型由来の配列であれば、最後の集計の際には同じ行にまとめるべきです。ラン ID とタグ ID の両方が異なっていると、その操作が面倒になります。タグ ID よりもラン ID の方が普通は少なく後から操作しやすいので、タグ ID (とプライマー ID) をここで統一しておきます。これ以降の処理は 1 ランでシーケンスした場合と同様になりますが、最後の集計表作成は以下のコマンドで行なって下さい。異なるラン ID の部分を一つのラン ID で置き換えてしまうことにより、ラン ID とタグ ID とプライマー ID が一致するようになり、単一の行にまとめて集計が行われます。

```
$ clsumclass \
--runname=HOGHEHOGE \
入力ファイル \
出力ファイル
```

2.2 Removal of noisy and/or chimeric sequences based on sequence dereplication

Claident では、配列のクラスターリング結果から読み間違いの多いノイジーな配列を推定して除去することができます。方法は CD-HIT-OTU (Li *et al.*, 2012) とほとんど同じです。また、UCHIME (Edgar *et al.*, 2011) を呼び出して結果を受け取ることで、キメラ配列の検出・除去も可能です。

以下のコマンドでノイジー配列とキメラ配列の検出・除去が行えます。なお、多数のファイルを一度に与えることができますが、できるだけ 1 ラン分のファイルを与えて下さい。これは、ランごとに塩基配列決定の質にばらつきがあるためです。入力ファイルのラン ID は統一されている必要があります。MiSeq など 1 ラン分の配列が非常に多い場合は、1 サンプルごとにこの処理を行って下さい。

```
$ clcleanseq \
uchime --minh 0.1 --mindiv 0.8 end \
--detectmode=N+C \
--pnoisycluster=0.5 \
--numthreads=使用する CPU 数 \
入力ファイル 1 \
中略 \
入力ファイル N \
```

出力フォルダ

オプションのうち、**uchime** から **end** の間は UCHIME の実行オプションです。--**minh** と --**mindiv** はキメラ配列と判定する感度に関するオプションです。--**minh** は値が小さいほど感度が高くなります。UCHIME の作者は 0.1 以上 5 以下の範囲にするよう推奨しています。デフォルト値は 0.1 です。--**mindiv** の指定値が 0.8 の場合、0.8% 以下の違いしかない配列間はキメラが形成されていても無視されます。最終的に 97% 一致規準で配列をクラスタリングするとしたら、2% 程度の違いしかない配列間のキメラはクラスタリング・コンセンサス配列作成によって潰されるため、この値を大きくすることで計算が速くなります。また、--**pnoisycluster** はノイジー配列と判定する感度に関するオプションです。0 以上 1 以下で指定して下さい。大きいほど感度が上がります。デフォルト値・推奨値は 0.5 です。クラスタリングを最終的に 97% 以下の閾値で行うのであればこれでいいと思いますが、99% 前後でクラスタリングする場合、--**pnoisycluster** は 0.9 くらいにするといいと思います。ただし、レア OTU を捨ててしまう可能性が高くなるので、1 サンプル当たりの配列数を多めに読んでおく必要があります。

各入力ファイルは同一の鋳型を同一のチューブ内で PCR した配列群を含んでいる必要があります。Claident がキメラ検出に利用している UCHIME では、元配列数が多いコンセンサス配列ほど PCR 当初から存在した=キメラではないと仮定し、元配列数が少ないコンセンサス配列が、より元配列数の多いコンセンサス配列の部分配列の組み合わせで説明できるかどうかを検討します。したがって、UCHIME に「同一の鋳型を同一のチューブ内で PCR した配列群内でのコンセンサス配列の元配列数」を正確に伝える必要があるため、1 ファイルは同一の鋳型を同一のチューブ内で PCR した配列群でなくてはならないのです。これは 1 ラン分のファイルを一度に与えることよりもはるかに優先度が高いので、同一の鋳型を同一のチューブ内で PCR したサンプルを複数回のランで塩基配列決定した場合などは複数ラン分のデータを入力しても構いません。

UCHIME によるキメラ判定はそれぞれのサンプルごとに行われ、複数のサンプルで観測されている完全一致配列の場合、どれか 1 つのサンプルでキメラと判定されていれば、他のサンプルでもキメラとみなして除去します。ただし、下記のオプションでこの動作を変更することができます。両オプションの両方を満たせばキメラと判定されます。

--**minnpositive** 整数値で指定。この値以上の配列 (サンプルではない) がキメラと判定されていれば、全サンプルでキメラと見なす。デフォルト値は 1。
 --**minppositive** 小数値で指定。キメラと判定された配列数/完全一致した全配列数がこの値以上であれば、全サンプルでキメラと見なす。デフォルト値は 0。

出力フォルダには、以下のファイル群が保存されます。

ラン ID__タグ ID__プライマー ID.chimeraremoved.dereplicated.fastq.gz キメラ配列除去し完全一致配列をまとめた配列
 ラン ID__タグ ID__プライマー ID.chimeraremoved.fastq キメラ配列除去した配列
 ラン ID__タグ ID__プライマー ID.denoised.dereplicated.fastq.gz ノイジー配列除去し完全一致配列をまとめた配列
 ラン ID__タグ ID__プライマー ID.denoised.fastq.gz ノイジー配列除去した配列
 ラン ID__タグ ID__プライマー ID.cleaned.dereplicated.fastq.gz キメラ配列とノイジー配列を除去し完全一致配列をまとめた配列
 ラン ID__タグ ID__プライマー ID.cleaned.fastq.gz キメラ配列とノイジー配列を除去した配列
 ラン ID__タグ ID__プライマー ID.chimericreads.txt.gz キメラ配列と判定された配列のリスト
 ラン ID__タグ ID__プライマー ID.noisyreads.txt.gz ノイジー配列と判定された配列のリスト
 ラン ID__タグ ID__プライマー ID.singletons.txt.gz 完全一致配列をまとめたときにシングルトンになった配列のリスト
 ラン ID__タグ ID__プライマー ID.uchime.txt.gz UCHIME によるキメラ判定結果のログ
 ラン ID__タグ ID__プライマー ID.parameter.txt まとまった完全一致配列をノイジーと判定するのに使用した所属配列数の下限値

キメラ配列除去を無効化してノイジー配列除去のみを行う場合は--**detectmode=noise**、逆にノイジー配列除去を無効化してキメラ配列除去のみを行う場合は--**detectmode=chimera** オプションを付けて実行して下さい。ただし、--**detectmode=noise** で処理した結果を--**detectmode=chimera** で処理しても、一度にノイジー配列・キメラ配列除

去を行った場合 (`--detectmode=N+C` を指定した場合またはデフォルト設定) と結果は一致しませんのでご注意ください。処理する場合は同時に行う必要があります。これは、キメラ配列もノイジー配列と判定するための情報を持っており、同時処理の場合にはそれを利用するためです。

2.3 Sequence clustering within sample

以下のコマンドで、同一の鋳型を同一のチューブ内で PCR した配列群のクラスタリングを行います。前節の結果得られる、完全一致配列をまとめた結果 (ラン ID__タグ ID__プライマー ID.cleaned.dereplicated.fastq.gz) を入力します。

```
$ clclasseq \  
--minident=0.99 \  
--strand=plus \  
--numthreads=使用する CPU 数 \  
入力ファイル 1 \  
中略 \  
入力ファイル N \  
出力フォルダ
```

前節の出力ファイルに対してこの処理を適用する場合、以下のようにコマンドを打てばいいでしょう。

```
$ clclasseq \  
--minident=0.99 \  
--strand=plus \  
--numthreads=使用する CPU 数 \  
"前節の出力フォルダ/*.cleaned.dereplicated.fastq.gz" \  
出力フォルダ
```

オプションのうち、`--minident=0.99` は 99% 以上一致する配列をまとめる設定です。最終的に 97% で配列をまとめる場合もこの時点では 99% にしておくのがおすすめです。最終的に 99% でまとめる場合にはここでは 99.5% にしておきます。`--strand=plus` は同ーストランドの比較だけでクラスタリングを行うオプションです。片側からしか読んでいないので、両ストランドの比較をする必要はなく、比較するとかえって異常なクラスタリングをしかねません。両側から読んでいる場合も、連結した後はストランドは揃っているのもこれでいいはず。出力フォルダ内の `*.assembled.fastq.gz` がクラスタリング後の配列です。

2.4 Sequence clustering among sample

サンプル内クラスタリングデータを用いてサンプル間でのクラスタリングを行い、全体のクラスタリング結果を得るには以下のコマンドを実行します。1 サンプルしかない場合は、上記のサンプル内配列クラスタリングを飛ばしてこの処理を行って下さい。また、多サンプルの場合でも、サンプル内配列クラスタリングの処理を飛ばしてこの処理を行っても構いませんが、サンプル内配列クラスタリングを行ってから処理した方が高速だと思います。

```
$ clclassclass \
--minident=0.99 \
--strand=plus \
--numthreads=使用する CPU 数 \
入力ファイル 1 \
中略 \
入力ファイル N \
出力フォルダ 1
$ clreclassclass \
--minident=0.97 \
--strand=plus \
--numthreads=使用する CPU 数 \
出力フォルダ 1 \
出力フォルダ 2
```

ここでの入力ファイルは前節の出力ファイル*.assembled.fastq.gz です。したがって、実際には以下のように入力することになるでしょう。

```
$ clclassclass \
--minident=0.99 \
--strand=plus \
--numthreads=使用する CPU 数 \
"前節の出力フォルダ/*.assembled.fastq.gz" \
出力フォルダ 1
$ clreclassclass \
--minident=0.97 \
--strand=plus \
--numthreads=使用する CPU 数 \
出力フォルダ 1 \
出力フォルダ 2
```

1 つ目の **clclassclass** でサンプル内配列クラスタリングでまとめられた配列において、99% 以上一致する配列をまとめます。2 つ目の **clreclassclass** で、97% 以上一致する配列をまとめ直します。こうすることで、より高速に、より正確にクラスタリングが行われます。**clreclassclass** をせずに **clclassclass** に 97% でまとめるよう指示しても構いませんが、アルゴリズム上わずかにまとめ漏れが発生します。

出力フォルダには、以下のファイルが保存されています。

assembled.contigmembers.gz どの生配列がどの OTU に割り当てられたかを記録しているファイル
assembled.fastq.gz クラスタリングされた配列
assembled.fasta クラスタリングされた配列

2.5 Summarizing result of clustering and postprocessing

2.5.1 Generating OTU x Sample table

以下のコマンドで、クラスタリング結果から各サンプルにおける各 OTU に割り当てられた生配列数を表にすることができます。


```
$ clsumclass \
--output=Matrix \
入力ファイル \
出力ファイル
```

ここでの入力ファイルは、クラスタリング結果の出力フォルダに保存されている、`assembled.contigmembers.gz` というファイルです。出力されるファイルは表 2.1 のようなタブ区切りのテキストファイルで、Excel などの表計算ソフトや R で読み込むことができます。ただし、表計算ソフトは巨大な行列の全てを読み込むことができない可能性がありますので注意して下さい。このファイルに基づいて、群集生態学的な解析を行うことができます。

samplename	amongsampleotu.1	amongsampleotu.2	amongsampleotu.3
sampleA	2371	0	0
sampleB	0	1518	0
sampleC	1398	0	0
sampleD	0	1436	0
sampleE	0	0	1360
sampleF	0	0	977

Table 2.1 集計表の例 — 数値は各サンプルにおける各 OTU の観測配列数です。

2.5.2 Eliminating OTUs and/or samples from table

`clsumclass` では、1 回だけ出現した配列も含めて全て出力されますが、以下のコマンドを用いてサンプルや OTU を選別することができます。入力ファイルは `clsumclass` の出力ファイルです。

```
$ clfiltersum \
オプション \
入力ファイル \
出力ファイル
```

使用できるオプションは以下の通りです。

- `--minnseqotu` 整数値で指定。この値に基づいて「割り当てられた生配列数がどのサンプルでも指定値未満の OTU」を除去します。
- `--minpseqotu` 小数値で指定。この値に基づいて「割り当てられた生配列数/サンプルの生配列総数がどのサンプルでも指定値未満の OTU」を除去します。
- `--minntotalseqotu` 整数値で指定。この値に基づいて「割り当てられた生配列総数が指定値未満の OTU」を除去します。
- `--minnseqsample` 整数値で指定。この値に基づいて「どの OTU の生配列数も指定値未満のサンプル」を除去します。
- `--minpseqsample` 小数値で指定。この値に基づいて「そのサンプルのその OTU に割り当てられた生配列数/OTU の生配列総数がどの OTU でも指定値未満のサンプル」を除去します。
- `--minntotalseqsample` 整数値で指定。この値に基づいて「生配列総数が指定値未満のサンプル」を除去します。
- `--otu` 残したい OTU 名をカンマで区切って指定します。
- `--negativeotu` 除去したい OTU 名をカンマで区切って指定します。
- `--otulist` ファイル名を指定。1 行に 1 つの OTU 名を記したテキストファイルとして残したい OTU 名を指定します。
- `--negativeotulist` ファイル名を指定。1 行に 1 つの OTU 名を記したテキストファイルとして除去したい OTU 名を指定します。
- `--otuseq` ファイル名を指定。FASTA 形式の配列ファイルとして残したい OTU 名を指定します。
- `--negativeotuseq` ファイル名を指定。FASTA 形式の配列ファイルとして除去したい OTU 名を指定します。
- `--sample` 残したいサンプル名をカンマで区切って指定します。
- `--negativesample` 除去したいサンプル名をカンマで区切って指定します。
- `--samplelist` ファイル名を指定。1 行に 1 つのサンプル名を記したテキストファイルとして残したいサンプル名を指定します。

--negativesamplelist ファイル名を指定。1 行に 1 つのサンプル名を記したテキストファイルとして除去したいサンプル名を指定します。

--replicatelist ファイル名を指定。PCR レプリケート関係にあるサンプルをタブ区切りで同一行に記述する。出力で 1 つのサンプルに統合される。

--runname ラン ID を指定。集計表中のサンプル名に含まれるラン ID が全てこれに置換される。同一サンプル名になったサンプルは統合される。

後述するノイジー配列・キメラ配列の除去や他領域配列の除去を行った上で、残った配列を用いて **--otuseq** などのオプションを用いた OTU の選別を行う場合は、他の選別より先に行ってください。

2.5.3 Removal of noisy and/or chimeric sequences based on PCR replication method

PCR レプリケートを用意してノイジー配列・キメラ配列を検出・除去するには、まず同一鋳型 DNA 由来の PCR レプリケートはどのサンプルとどのサンプルなのかを記したファイルを用意する必要があります。以下のような内容で用意して下さい。

```
sample1 sample2 sample3
sample4 sample5
sample6 sample7
```

同一の行内にタブ区切りで書かれたサンプルが、同一鋳型由来の PCR レプリケートであることを表しています。異なる鋳型由来のサンプルは別の行に記述して下さい。PCR レプリケートは 3 つ以上あっても構いません。また、レプリケート数が鋳型ごとに異なっても構いません。

```
$ clfilterseq \
--contigmembers=*.contigmembers.gz \
--replicatelist=PCR レプリケート指示ファイル \
入力ファイル \
出力ファイル
```

レプリケートが 3 つ以上であっても、全てのレプリケートで検出された配列のみがノイジーでもキメラでもないと判断されます。また、1 つの OTU が複数の鋳型から検出されていた場合、どれか 1 つの鋳型でノイジーまたはキメラと判定されていれば、他の鋳型でもノイジーまたはキメラとみなして除去します。これらの設定は以下のオプションで変更することができます。

--minnreplicate 整数値で指定。この値以上のレプリケートで観測された OTU はそのサンプルではノイジーでもキメラでもないと見なす。デフォルト値は 2。

--minpreuplicate 小数値で指定。OTU が観測されたレプリケート数/サンプルのレプリケート数がこの値以上であれば、そのサンプルではノイジーでもキメラでもないと見なす。デフォルト値は 1。

--minnpositive 整数値で指定。この値以上の配列 (サンプルではない) がノイジーまたはキメラと判定されていれば、全サンプルでノイジーまたはキメラと見なす。デフォルト値は 1。

--minppositive 小数値で指定。ノイジーまたはキメラと判定された配列数/OTU の全配列数がこの値以上であれば、全サンプルでノイジーまたはキメラと見なす。デフォルト値は 0。

--runname ラン ID を指定。サンプル名に含まれるラン ID が全てこれに置換される。同一サンプル名になったサンプルは統合される。

--minnreplicate と **--minpreuplicate** の値に基づく 1 サンプル内での判定と、**--minnpositive** と **--minppositive** の値に基づく全サンプルにまたがる判定が 2 段階で行われていることに注意して下さい。前者では 2 つの条件を両方満たすとノイジーでもキメラでもないと判定され、後者では 2 つの条件を両方満たすとノイジーまたはキメラであると判定されます。また、同一の OTU が、あるサンプルではノイジーまたはキメラだが、別のサンプルではそうでない、

といった判定は行えません。最終的な判定は全サンプル共通です。また、PCR レプリケート指示ファイルに記されていないサンプルでは判定が行われず、全ての OTU がそのサンプルではノイズでもキメラでもないと仮定されます。

ここで作成した FASTA 配列ファイルを第 2.5.2 節の `clfilterseq` の `--otuseq` オプションに指定して集計表を加工することで、この配列ファイルから除去された OTU を集計表からも除去することができます。また、`--replicatelist` オプションと `--runname` オプションを `clfilterseq` と同様に指定すれば、PCR レプリケート関係にあるサンプルは 1 つのサンプルに統合されます (セルの値は合算されます)。

2.5.4 Removal of rarely observed OTU sequences

クラスタリングの出力フォルダには、`assembled.fasta` というファイルとして、最終的に得られたコンセンサス配列が保存されています。これを次章の DNA バーコーディングに用いますが、全サンプルを通して少量しか出現しないものまで含めると数が多すぎることがあります。そこで、以下のようにして 5 配列以上出現するものだけを抽出することができます。時間がかかってもいいのであれば、`assembled.fasta` をそのまま使っていただいても構いません。配列数の閾値を適当に変えてもいいでしょう。この閾値は慎重に決定して下さい。ノイズ配列・キメラ配列の除去を行っている場合には必要ないかもしれません。

```
$ clfilterseq \  
--contigmembers=*.contigmembers.gz \  
--minnseq=5 \  
入力ファイル \  
出力ファイル
```

ここで作成した FASTA 配列ファイルを第 2.5.2 節の `clfilterseq` の `--otuseq` オプションに指定して集計表を加工することで、この配列ファイルから除去された OTU を集計表からも除去することができます。

2.5.5 Dividing sequences by conserved region

ここまでの配列が複数領域にまたがっている場合 (たとえば ITS1・5.8S rRNA・ITS2 といった風に)、領域ごとに分割した方が良くかもしれません。領域間かその付近に保存的な領域があれば、それを境界の目印として領域を分割できます。特にユニバーサルプライマー配列などを使うといいでしょう。これは以下のコマンドで行うことができます。

```
$ cldivseq \  
--query=保存領域の配列 \  
--border=start \  
入力ファイル \  
出力ファイル 1 \  
出力ファイル 2
```

このコマンドでは、指定された配列を Needleman-Wunsch アルゴリズムを用いたアライメントによって 15% まで不一致を許して探します。該当する部位が見つかったら、その左端を境界として配列を 2 つに分け、それぞれ出力ファイル 1 と出力ファイル 2 に保存します。該当する部位が見つからなかった場合は、そのまま出力ファイル 1 だけに保存されます。信頼度配列ファイルがある場合には、こちらも塩基配列に合わせて分割されます。なお、指定する配列は対象配列と同一のストランドにしておいて下さい。ストランドが異なる場合は、`--reversecomplement` オプションを付

加して実行することで、`--query` に指定された配列の逆相補配列をアライメントに用いるようになります。

`--border` オプションを `start` ではなく `end` にすることで、検索に一致した部位の右端を境界として分割することができます。検索した配列が見つからなかった場合は出力ファイル 1 のみに保存され、出力ファイル 2 には配列が保存されませんが、これでは都合が悪い場合は `--makedummy` オプションを付加することでダミーの塩基配列が出力ファイル 2 に保存されます。分割した各領域の配列で別個に宿主生物を同定した上で同定結果を統合する場合など、配列間に対応がとれていないと困る場合にご利用下さい。ここでは 2 分割の例を挙げましたが、分割後にさらに分割を繰り返すことで 3 分割や 4 分割も可能です。

2.5.6 Extracting ITS or SSU rRNA sequences by ITSx or Metaxa

ITSx は核 ITS1・ITS2 領域を認識してその部位だけを取り出すことができるプログラムです (Bengtsson-Palme *et al.*, 2013)。多くの分類群で ITS は非常に変異に富んでいますが、ずっと保守的な SSU・LSU rRNA と隣接しているため、これらの配列が残っていると、その配列の影響で ITS では明確に区別できる遠縁な生物の SSU・LSU 配列が BLAST 検索時にヒットしてしまい、うまく同定できなくなってしまうことがあります。ITSx で ITS 領域のみを切り出して同定することでこのような問題に対処できます。

Metaxa は ITSx と同様に SSU (12S/16S/18S) rRNA を区別してそれぞれの部位だけを取り出すことができるプログラムです (Bengtsson *et al.*, 2011)。SSU rRNA は真核生物のメタバーコーディング・DNA バーコーディングに広く利用されている領域ですが、核だけでなくミトコンドリアや葉緑体、混入した細菌にも含まれているため、核の SSU を標的とする PCR を行った場合にも多少はミトコンドリアや細菌の SSU rRNA が混入してしまいます。そのような配列は群集生態学的解析の障害となるため、このプログラムで前もって除去しておくとの後の解析が楽になります。

いずれのプログラムもインストール方法、使用方法をここでは説明しません。それぞれの Web サイトとマニュアルをご覧ください。これらのプログラムで作成した FASTA 配列ファイルを第 2.5.2 節の `clfiltersum` の `--otuseq` オプションに指定して集計表を加工することで、この配列ファイルに含まれている OTU だけの集計表を作成することができます。

2.5.7 Removal of non-target sequences

ITSx や Metaxa は ITS と SSU rRNA 領域にしか使用できませんが、もっと汎用的な非ターゲット領域の探索方法があります。一つは BLAST などで検索して配列がどの遺伝子のものかを判定してくれるプログラムにかけることで、もう一つは多重整列プログラムを利用することです。ClustalW2・ClustalX2 (Larkin *et al.*, 2007) や MAFFT (Katoh & Standley, 2013) には、多重整列した結果を系統的に近いものが近くなるように並び替えて出力することができます。この機能を使って並び替えた上で、多重整列の閲覧ソフトで表示して目で確認すればどれが非ターゲット領域かすぐわかります。非ターゲット領域の OTU を除去した FASTA 配列ファイルを第 2.5.2 節の `clfiltersum` の `--otuseq` オプションに指定して集計表を加工することで、この配列ファイルに含まれている OTU だけの集計表を作成することができます。

2.6 Submission of sequence data to DRA

クラスタリング後の配列は従来通り DDBJ などに登録すればいいですが、新型シーケンサーの生データは DDBJ Sequence Read Archive (DRA) というところへ登録することになっています (サンガー法シーケンサーの生データは Trace Archive)。本書にあるように複数サンプルからの DNA にタグを付けてシーケンスした場合、サンプルごとのデータファイルに分割して登録することを求められます。これは `clsplitseq` によって分割した FASTQ を使えば問題ありません。

また、サンプルについての情報 (メタデータ) を記述した XML ファイルを作成しなくてはなりません。データ登録を受け付けている側でも、XML 作成を補助するツールが用意されていますが大量サンプルを扱っている場合は用意するのは大変です。そこで、DRA 登録用 XML 作成を補助する `clmaketsv` および `clmakexml` というコマンドを用意しています。以下でこれらの使用法を説明します。

DRA への登録には、DDBJ が管理する D-way でのユーザーアカウント作成と公開鍵の登録が必要です。詳しくは

<http://trace.ddbj.nig.ac.jp/dra/submission.shtml>

をご覧ください。また、登録の際、Submission、Study、Experiment、Sample、Run の概念と対応関係を把握している必要があります。これは、

<http://trace.ddbj.nig.ac.jp/dra/documentation.shtml>

をよく読んで予め理解しておいて下さい。

DRA では、Study は BioProject という研究プロジェクトデータベースに登録して、それを参照することになります。BioProject への登録は

<http://trace.ddbj.nig.ac.jp/bioproject/submission.html>

を参照して予め済ませておいて下さい。

さらに、Sample は BioSample という研究サンプルデータベースに登録して、それを参照します。BioSample への登録は

<http://trace.ddbj.nig.ac.jp/biosample/submission.html>

を参照して予め済ませておいて下さい。なお、メタゲノムを鋳型としてユニバーサルプライマーで増幅を行ったシーケンスサンプルの場合、MIMarks-Survey を MIMixS のタイプとして指定します。サンプルの採集された高度、水深、温度、湿度、pH など、様々な情報を付加しておくことができます。後世のためにも、面倒がらずにできるだけ多くの情報を付加しておいて下さい。

2.6.1 Generating tab-delimited text

DRA に登録する XML には、登録する FASTQ ファイルごとにそのファイルに含まれている配列の情報を記述します。これは非常に複雑で大変なため、一旦単純なタブ区切りテキストに書き出した上で、Excel などを用いてそれを編集し、編集後のタブ区切りテキストから XML を作成します。タブ区切りテキストを生成するコマンドは `clmaketsv` です。以下のように使って下さい。

```
$ clmaketsv \  
入力ファイル 1 \  
中略 \  
入力ファイル N \  
出力ファイル
```

入力ファイルは登録するものを指定して下さい。ワイルドカードも使えます。タブ区切りテキストができれば、表計算ソフトなどに読み込ませて、各セルを埋めていって下さい。各セル内では、[Hogehoge,Fugafuga] とあった場合は Hogehoge または Fugafuga のどちらかを選んで括弧と選ばなかったものを消して下さい。<Fill in this cell>とあった場合は、<>内の指示に従ってセルを埋めて下さい (括弧は残さない)。その他のセルは適当に自分で考えて何とかして下さい。

2.6.2 Generating DRA XML from tab-delimited text

タブ区切りテキストの編集が終わったら、タブ区切りテキストとして保存した上で以下のように `clmakexml` を実行して下さい。登録に必要な XML ファイル群が生成されます。なお、タブ区切りテキストは一度に複数指定することも可能です。ただし、各ファイルの最初の 3 行に記述する内容は、最初のファイルのものしか利用されません。2 ファイル目以降の最初の 3 行は無視されます。

```
$ clmakexml \  
タブ区切りテキストファイル \  
DRA から割り振られた submission-ID
```

DRA では、「Create new submission(s)」によって submission-ID が割り振られます。ユーザー ID-0001 などとなっているはずですが、DRA から割り振られた submission-ID にはこれを指定します。すると、submission-ID.*.xml という名前のファイルが 3 つ生成されます。これを DRA の XML Upload 機能を利用して送信します。その他もろもろの処理が終わると、登録してあるメールアドレスに割り振られたアクセス番号が送られてきます。論文にはそれを書いておけばいいでしょう。

Chapter 3

Identification of host organisms of sequences by DNA barcoding

DNA バーコーディングは、近年非常に多くの分野で応用が進められている、生物の同定方法です。しかし、既知配列データベースが不十分な上、それを前提とした同定アルゴリズムが欠けていました。そこで、筆者は「問い合わせ配列と最近隣配列との間の変異量<分類群内の最大変異量」という規準を考案し、これを実現するアルゴリズム QCauto 法 (Tanabe & Toju, 2013) を Clident に実装しました。ただし、その配列のホストの可能性のある全生物が記載済みで、しかもそのバーコード配列もデータベースに登録済みである場合には、そのような難しいことを考えずに最近隣配列と同種とみなせばよいでしょう。そちらの方法も Clident には実装してあります。以下のコマンドは全てターミナルかコンソールで実行して下さい。基本的なターミナルの使い方の知識は持っているものとして話を進めます。

3.1 Retrieving neighborhood sequences by BLAST search

以下のコマンドで、「問い合わせ配列と最近隣配列との間の変異量<分類群内の最大変異量」という条件を満たすために検討すべき近隣既知配列群 (の GenBank ID) を取得できます。

```
$ clidentseq \  
--blastdb=overall_genus \  
--numthreads=使用する CPU 数 \  
入力ファイル \  
出力ファイル
```

入力ファイルには FASTA 形式の塩基配列ファイルを指定します。--blastdb オプションでは、BLAST 検索に用いるデータベース名を指定します。筆者が用意しているデータベースは以下の通りです。

animals_COX1_genus 動物 (Metazoa) の mtDNA COX1 配列で属以下の情報があるもの
animals_COX1_species 同上だが種以下の情報があるもの
animals_mt_genus 動物 (Metazoa) の mtDNA 配列で属以下の情報があるもの
animals_mt_species 同上だが種以下の情報があるもの
eukaryota_LSU_genus 真核生物の LSU (28S) rRNA 配列で属以下の情報があるもの
eukaryota_LSU_species 同上だが種以下の情報があるもの
eukaryota_SSU_genus 真核生物の SSU (18S) rRNA 配列で属以下の情報があるもの
eukaryota_SSU_species 同上だが種以下の情報があるもの
fungi_ITS_genus 真菌の ITS 配列で属以下の情報があるもの

fungi.ITS.species 同上だが種以下の情報があるもの
overall.class NCBI nt の中で綱以下の情報があるもの
overall.order 同上だが目以下の情報があるもの
overall.family 同上だが科以下の情報があるもの
overall.genus 同上だが属以下の情報があるもの
overall.species 同上だが種以下の情報があるもの
plants.matK.genus 緑色植物の cpDNA *matK* 配列で属以下の情報があるもの
plants.matK.species 同上だが種以下の情報があるもの
plants.rbcL.genus 緑色植物の cpDNA *rbcL* 配列で属以下の情報があるもの
plants.rbcL.species 同上だが種以下の情報があるもの
plants.trnH-psbA.genus 緑色植物の cpDNA *trnH-psbA* 配列で属以下の情報があるもの
plants.trnH-psbA.species 同上だが種以下の情報があるもの
prokaryota.16S.genus 原核生物の 16S rRNA 配列で属以下の情報があるもの
prokaryota.16S.species 同上だが種以下の情報があるもの
semiall.class **overall.class** から脊椎動物と *Caenorhabditis* と *Drosophila* を除いたもの
semiall.order **overall.order** から脊椎動物と *Caenorhabditis* と *Drosophila* を除いたもの
semiall.family **overall.family** から脊椎動物と *Caenorhabditis* と *Drosophila* を除いたもの
semiall.genus **overall.genus** から脊椎動物と *Caenorhabditis* と *Drosophila* を除いたもの
semiall.species **overall.species** から脊椎動物と *Caenorhabditis* と *Drosophila* を除いたもの

overall.* と **semiall.*** は非常に巨大なため、多くのメモリを必要とします。**overall.*** で 20GB 弱、**semiall.*** で 10GB 弱です。これらは万能なのでどんな配列にも使えますし、想定外の分類群や想定外の遺伝子座の配列の混入があってもとんでもない誤同定を避けることができます。しかしメモリ要求量が多いため、多くの CPU を使用しようとして BLAST をたくさん起動すると、起動した BLAST 数だけさらにメモリを多く必要とします。つまり、**overall.*** を使った BLAST 検索で 4 つの CPU を使うには、80GB ものメモリが必要になります。ただし、BLAST ではメモリ不足で HDD や SSD にスワップしてしまってもそれほど遅くならないので、64GB 程度のメモリと大きめのスワップ領域があれば実行可能です。また、`--numthreads=4 --hyperthreads=2` というオプションを使うと、BLAST を 2 つ起動して、2 つの BLAST それぞれに 2 つの CPU を使わせます。この場合は必要なメモリは 40GB 程度に減少します。スワップがあれば 32GB でも実行可能だと思います。**overall.genus** と **semiall.genus** では属以下の情報がある、つまりそれなりに分類群の情報が信頼できそうな配列を選別していますが、あまりにマイナーな分類群だとそれでは近隣配列が見つからないことがあります。その場合は、綱以下まで同定された配列のデータベース **overall.class** か **semiall.class** を使ってみて下さい。その他のデータベースは、**overall.*** や **semiall.*** よりもメモリ占有量がずっと少ないため、メモリの少ないコンピュータでも動作します。

問い合わせ配列のホストの可能性のある全生物が記載済みで、しかもそのバーコード配列もデータベースに登録済みである場合には、以下のようにして類似度 99% 以上の上位 1 位 (1 位タイを含む) を近隣配列とすればいいでしょう。

```

$ clidentseq \
--method=1,99% \
--blastdb=overall.genus \
--numthreads=使用する CPU 数 \
入力ファイル \
出力ファイル

```

3.2 Taxonomic assignment based on neighborhood sequences

以下のコマンドで、近隣既知配列群の所属分類群が同一になるまで分類階層を上げていくことで配列を同定します。これは lowest common ancestor (LCA) algorithm と呼ばれています (Huson *et al.*, 2007)。


```
$ classigntax \
--taxdb=overall_genus \
入力ファイル \
出力ファイル
```

入力ファイルには前節で作成した **clidentseq** の出力ファイルを指定して下さい。--**taxdb** は既知配列の所属分類群データベース指定オプションです。BLAST データベースと同じ名前のものがインストールされていますので、それを指定して下さい。

classigntax のデフォルト設定では、少なくとも 2 本以上の近隣既知配列がないと同定することができません。前節後半のように **clidentseq** の実行時に --**method=1,99%** を指定していると、近隣既知配列が 1 本しかないために一つも同定できないことになります。この場合は以下のようにして必要な近隣既知配列数を 1 にします。

```
$ classigntax \
--taxdb=overall_genus \
--minnsupporter=1 \
入力ファイル \
出力ファイル
```

出力ファイルは表 3.1 のような形のタブ区切りのテキストになっています。

query	phylum	genus	species
seqA	Ascomycota	<i>Chloridium</i>	<i>Chloridium virescens</i>
seqB	Ascomycota	<i>Chloridium</i>	<i>Chloridium virescens</i>
seqC	Ascomycota	<i>Chloridium</i>	<i>Chloridium virescens</i>
seqD	Basidiomycota	<i>Amanita</i>	<i>Amanita fuliginea</i>
seqE	Basidiomycota	<i>Coltriciella</i>	<i>Coltriciella dependens</i>
seqF	Basidiomycota	<i>Filobasidium</i>	<i>Filobasidium uniguttulatum</i>
seqG	Basidiomycota	<i>Laccaria</i>	<i>Laccaria bicolor</i>
seqH	Basidiomycota	<i>Lactarius</i>	<i>Lactarius quietus</i>
seqI	Basidiomycota	<i>Russula</i>	<i>Russula densifolia</i>
seqJ	Basidiomycota	<i>Russula</i>	<i>Russula densifolia</i>
seqK	Basidiomycota	<i>Russula</i>	<i>Russula densifolia</i>
seqL	Basidiomycota	<i>Russula</i>	<i>Russula vesca</i>
seqM	Basidiomycota	<i>Agaricus</i>	
seqN	Basidiomycota	<i>Amanita</i>	
seqO	Basidiomycota	<i>Amanita</i>	
seqP	Ascomycota	<i>Bisporrella</i>	
seqQ	Ascomycota	<i>Capronia</i>	
seqR	Ascomycota	<i>Capronia</i>	
seqS	Ascomycota	<i>Cenococcum</i>	

Table 3.1 同定結果の例 — 空欄は **unidentified**、つまり不明ということです。横幅を抑えるためにいくつかの分類階層は省いてあります。

classigntax のデフォルト設定では、全ての近隣既知配列の所属分類群が同一になるまで分類階層を上げていきますので、誤同定された配列が混ざっていたりした場合には下位の分類階層はほとんど不明になってしまいます。これは「ほぼ正しいと思われる」結果を得ることをデフォルト設定では優先しているためですが、近隣既知配列が多い場合、多少不一致を許容して誤同定の可能性が増してでもできるだけ下位の分類階層まで同定したい場合があります。そのよ

うな場合には以下のように`--maxpopposer`と`--minsortatio`を指定します。

```
$ classigntax \
--taxdb=overall_genus \
--maxpopposer=0.05 \
--minsortatio=19 \
入力ファイル \
出力ファイル
```

`classigntax` は結果として採用する分類群に該当する配列を `supporter` 配列とし、該当しない配列を `opposer` 配列とします。`--maxpopposer` には `opposer` 配列の存在を許容する割合を百分率で指定します。`--minsortatio` は `supporter` 配列数と `opposer` 配列数の比の下限値を指定します。`supporter` 配列数と `opposer` 配列数がこれを下回るような同定結果は許容せず分類階層を上げていきます。2つのオプションが必要なのは、その分類階層の情報がない配列が存在し得るためです。そのような配列は `supporter` でも `opposer` でもないので、2つのオプションが必要になります。上記の例では、近隣既知配列中、5% までの `opposer` 配列を許容し、`opposer` 配列の 19 倍以上の `supporter` 配列を必要とします。

3.3 Integrating multiple identification results

例えば植物の *rbcL* と *matK* など、複数のバーコード領域を併用しなければ同定できないこともありますし、`overall_genus` と `overall_class` での同定結果のいいところ取りをしたい場合があります。また、厳密な LCA による同定結果と 5% まで `opposer` を許容する LCA の同定結果の組み合わせでいいところ取りしたい場合もあるでしょう。そのようなことが、複数の同定結果を統合することで可能になります。

植物の *rbcL* と *matK* を併用して同定する場合は、より深くまで同定できている方の結果を採用するのがいいでしょう。これは以下のコマンドでできます。

```
$ clmergeassign \
--priority=equal \
--preferlower \
rbcL での同定結果 \
matK での同定結果 \
出力ファイル
```

より保守的に考えるなら、以下のようにして、両方の同定結果が同一か、一方では未同定の場合だけ採用することもあります。ただし、上の階層で不一致だった場合には、同定結果が一方で未同定でも採用されません。

```
$ clmergeassign \
--priority=equal \
rbcL での同定結果 \
matK での同定結果 \
出力ファイル
```

trnH-psbA も併用して、最も深くまで同定できているものを採用する場合は以下のようにします。

```
$ clmergeassign \  
--priority=equal \  
--preferlower \  
rbcL での同定結果 \  
matK での同定結果 \  
trnH-psbA での同定結果 \  
出力ファイル
```

overall_genus の結果を優先的に採用し、**overall_genus** で同定できている分類階層までは一致しているが **overall_class** の方がより下位まで同定できている場合だけ採用するには、以下のようにします。

```
$ clmergeassign \  
--priority=descend \  
overall_genus での同定結果 \  
overall_class での同定結果 \  
出力ファイル
```

同様に、厳密な LCA による結果を優先し、厳密な LCA で同定できている分類階層まで一致しているが制約を緩めた LCA による結果の方がより下位まで同定できている場合にのみ採用するには、以下のようにします。

```
$ clmergeassign \  
--priority=descend \  
厳密な LCA での同定結果 \  
制約を緩めた LCA での同定結果 \  
出力ファイル
```


References

- Andrews, S. (2010) Software distributed by the author at <http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>.
- Bengtsson, J., Eriksson, K. M., Hartmann, M., Wang, Z., Shenoy, B. D., Grelet, G.-A., Abarenkov, K., Petri, A., Rosenblad, M. A., Nilsson, R. H. (2011) Metaxa: a software tool for automated detection and discrimination among ribosomal small subunit (12S/16S/18S) sequences of archaea, bacteria, eukaryotes, mitochondria, and chloroplasts in metagenomes and environmental sequencing datasets. *Antonie Van Leeuwenhoek*, **100**, 471–475.
- Bengtsson-Palme, J., Ryberg, M., Hartmann, M., Branco, S., Wang, Z., Godhe, A., De Wit, P., Sánchez-García, M., Ebersberger, I., de Sousa, F., Amend, A., Jumpponen, A., Unterseher, M., Kristiansson, E., Abarenkov, K., Bertrand, Y. J. K., Sanli, K., Eriksson, K. M., Vik, U., Veldre, V., Nilsson, R. H. (2013) Improved software detection and extraction of ITS1 and ITS2 from ribosomal ITS sequences of fungi and other eukaryotes for analysis of environmental sequencing data *Methods in Ecology and Evolution*, **4**, 914–919.
- Edgar, R. C., Haas, B. J., Clemente, J. C., Quince, C., Knight, R. (2011) UCHIME improves sensitivity and speed of chimera detection. *Bioinformatics*, **27**, 2194–2200.
- Fadrosh, D. W., Ma, B., Gajer, P., Sengamalay, N., Ott, S., Brotman, R. M., Ravel, J. (2014) An improved dual-indexing approach for multiplexed 16S rRNA gene sequencing on the Illumina MiSeq platform. *Microbiome*, **2**, 6.
- Hamady, M., Walker, J. J., Harris, J. K., Gold, N. J., Knight, R. (2008) Error-correcting barcoded primers for pyrosequencing hundreds of samples in multiplex. *Nature Methods*, **5**, 235–237.
- Huson, D. H., Auch, A. F., Qi, J., Schuster, S. C. (2007) MEGAN analysis of metagenomic data. *Genome Research*, **17**, 377–386.
- Illumina corporation (2013) 16S metagenomic sequencing library preparation.
- Katoh, K., Standley, D. M. (2013) MAFFT multiple sequence alignment software version 7: improvements in performance and usability. *Molecular Biology and Evolution*, **30**, 772–780.
- Kunin, V., Engelbrektson, A., Ochman, H., Hugenholtz, P. (2010) Wrinkles in the rare biosphere: pyrosequencing errors can lead to artificial inflation of diversity estimates. *Environmental Microbiology*, **12**, 118–123.
- Larkin, M. A., Blackshields, G., Brown, N. P., Chenna, R., McGettigan, P. A., McWilliam, H., Valentin, F., Wallace, I. M., Wilm, A., Lopez, R., Thompson, J. D., Gibson, T. J., Higgins, D. G. (2007) Clustal W and Clustal X version 2.0. *Bioinformatics*, **23**, 2947–2948.
- Li, W., Fu, L., Niu, B., Wu, S., Wooley, J. (2012) Ultrafast clustering algorithms for metagenomic sequence analysis. *Briefings in Bioinformatics*, **13**, 656–668.
- Nelson, M. C., Morrison, H. G., Benjamino, J., Grim, S. L., Graf, J. (2014) Analysis, optimization and verification of Illumina-generated 16S rRNA gene amplicon surveys. *PLoS One*, **9**, e94249.
- Schmieder, R., Edwards, R. (2011) Quality control and preprocessing of metagenomic datasets. *Bioinformatics*, **27**, 863–864.
- Stevens, J. L., Jackson, R. L., Olson, J. B. (2013) Slowing PCR ramp speed reduces chimera formation from environmental samples. *Journal of Microbiological Methods*, **93**, 203–205.

- Tanabe, A. S., Toju, H. (2013) Two new computational methods for universal DNA barcoding: a benchmark using barcode sequences of bacteria, archaea, animals, fungi, and land plants. *PLoS One*, **8**, e76910.
- Zhang, J., Kobert, K., Flouri, T., Stamatakis, A. (2014) PEAR: a fast and accurate Illumina Paired-End read mergeR. *Bioinformatics*, **30**, 614–620.

Appendix A

Installing optional programs

A.1 Installation of bcl2fastq

Illumina から提供されている、BCL 形式のベースコールデータから FASTQ を生成するプログラム bcl2fastq は MiSeq・HiSeq 用の v1.8.4 および NextSeq 500・HiSeq X 用の v2.15 が

http://support.illumina.com/downloads/bcl2fastq_conversion_software.html

からダウンロードできます。RedHat 系 Linux 用の .rpm で提供されていますので、Ubuntu を含む Debian 系では alien というプログラムで Debian 用の .deb に変換しインストールします。v1.8.4 のインストールをする場合は以下のようにコマンドを実行して下さい。

```
# alien のインストール
$ sudo apt-get install alien
# bcl2fastq v1.8.4 のダウンロード
$ wget \
ftp://webdata.webdata@ussd-ftp.illumina.com/Downloads/Software/bcl2fastq/bcl2fastq-1.8.4-Linux-x86_64.rpm
# .rpm から .deb への変換とインストール
$ sudo alien -i \
bcl2fastq-1.8.4-Linux-x86_64.rpm
# bcl2fastq に必要な追加パッケージのインストール
$ sudo apt-get install libxml-simple-perl xsltproc
```

v2.15 をインストールしたい場合は、前述の Web ページからダウンロードして上記と同様にインストールするか、ダウンロード URL とファイル名を変更して実行して下さい。なお、使用方法が v1.8.4 とは大幅に異なりますので、マニュアルをよく読んでご使用下さい。

Ubuntu・Xubuntu・Lubuntu の 14.04 LTS では、Perl のバージョンが新しすぎ、より古いバージョンを前提として書かれた bcl2fastq が動作しません。そこで以下のようにして問題の部分を書き換えます。

```
$ sudo perl -i -npe \
's/qw\(\ELAND_FASTQ_FILES_PER_PROCESS\)\(/"ELAND_FASTQ_FILES_PER_PROCESS"/'\ \
/usr/local/lib/bcl2fastq-1.8.4/perl/Casava/Alignment/Config.pm
$ sudo perl -i -npe \
's/qw\(\ELAND_GENOME\)\(/"ELAND_GENOME"/'\ \
/usr/local/lib/bcl2fastq-1.8.4/perl/Casava/Alignment/Config.pm
```

テキストエディタで

`/usr/local/lib/bcl2fastq-1.8.4/perl/Casava/Alignment/Config.pm`

の 747 行目と 751 行目を手動で書き換えても構いません。`qw(HOGEHOGE)` を `("HOGEHOGE")` に書き換えて下さい。

Appendix B

Useful terminal commands

B.1 Counting sequences

```
# FASTQ の場合
$ grep -P -c '^\\+$' inputfile
# GZIP 圧縮 FASTQ の場合
$ gzip -dc inputfile | grep -P -c '^\\+$'
# BZIP2 圧縮 FASTQ の場合
$ bzip2 -dc inputfile | grep -P -c '^\\+$'
# XZ 圧縮 FASTQ の場合
$ xz -dc inputfile | grep -P -c '^\\+$'
# FASTA の場合
$ grep -P -c '^>' inputfile
# GZIP 圧縮 FASTA の場合
$ gzip -dc inputfile | grep -P -c '^>'
# BZIP2 圧縮 FASTA の場合
$ bzip2 -dc inputfile | grep -P -c '^>'
# XZ 圧縮 FASTA の場合
$ xz -dc inputfile | grep -P -c '^>'
```

B.2 Viewing sequences

```
# 無圧縮ファイル内容を画面に出力
$ cat inputfile
# 無圧縮ファイル内容をスクロールしながら見る
$ less inputfile
# GZIP 圧縮ファイル内容を画面に出力
$ gzip -dc inputfile
# GZIP 圧縮ファイル内容をスクロールしながら見る
$ gzip -dc inputfile | less
# FASTQ の最初の 1 配列を画面に出力
$ head -n 4 inputfile
# GZIP 圧縮 FASTQ の最初の 1 配列を画面に出力
$ gzip -dc inputfile | head -n 4
# FASTQ の特定の配列を検索して画面に出力
$ grep -P -A 3 '^@配列名の正規表現' inputfile
# GZIP 圧縮 FASTQ の特定の配列を検索して画面に出力
```

```
$ gzip -dc inputfile | grep -P -A 3 '^@配列名の正規表現'
```

B.3 Compression and decompression

```
# フォルダ内の全.fastq.gz を展開
$ for f in *.fastq.gz
do gzip -d $f
done
# フォルダ内の全.fastq.gz を展開 (サブフォルダ含む)
$ for f in $(find . -name *.fastq.gz)
do gzip -d $f
done
# 4つのCPUを使用してフォルダ内の全.fastq.gz を展開
$ ls *.fastq.gz | xargs -L 1 -P 4 gzip -d
# 4つのCPUを使用してフォルダ内の全.fastq.gz を展開 (サブフォルダ含む)
$ find . -name *.fastq.gz | xargs -L 1 -P 4 gzip -d
# フォルダ内の全.fastq を圧縮
$ for f in *.fastq
do gzip $f
done
# フォルダ内の全.fastq を圧縮 (サブフォルダ含む)
$ for f in $(find . -name *.fastq)
do gzip $f
done
# 4つのCPUを使用してフォルダ内の全.fastq を圧縮
$ ls *.fastq | xargs -L 1 -P 4 gzip
# 4つのCPUを使用してフォルダ内の全.fastq を圧縮 (サブフォルダ含む)
$ find . -name *.fastq | xargs -L 1 -P 4 gzip
```

B.4 Sequence selection

```
# FASTQ から最初の 1 万配列を抽出
$ head -n 40000 inputfile > outputfile
# FASTQ から最後の 1 万配列を抽出
$ tail -n 40000 inputfile > outputfile
# GZIP 圧縮 FASTQ から最初の 1 万配列を抽出して GZIP 圧縮 FASTQ へ保存
$ gzip -dc inputfile | head -n 40000 | gzip -c > outputfile
# GZIP 圧縮 FASTQ から最後の 1 万配列を抽出して GZIP 圧縮 FASTQ へ保存
$ gzip -dc inputfile | tail -n 40000 | gzip -c > outputfile
# GZIP 圧縮 FASTQ の特定の配列を検索して GZIP 圧縮 FASTQ へ保存
$ gzip -dc inputfile | grep -P -A 3 '^@配列名の正規表現' | gzip -c > outputfile
```