

FASTQ ファイルのクオリティの統計を診るのに便利なソフトウェアに [FastQC](#) があります。

3 参照配列の準備

RNA-seq は転写産物のせいぜい数百 bp 程度の短い配列情報で、転写産物完全長 (~数 Kbp) の一部です。RNA-seq から遺伝子発現行列を得るには、ゲノムや完全長 cDNA の既知情報を参照配列として RNA-seq データを由来する転写産物別にまとめる作業が必要になります。そのためのレファレンスは、ゲノム・データベースや UniGene データベースなどから得ることができます。

テスト用の [参照配列データ](#) をダウンロードしてみましょう。これはイネの発現遺伝子のうち選抜した 344 個の配列を、FASTA 形式でまとめたものです。

```
>Os12t0274750-00 Hypothetical gene.
TCGCAGCCCGGGGCTTGTAGGCGATGAAGCTGATGAGCTGCACCTGCCTGACGTTGTGCG
...
>Os12t0115100-00 Nonspecific lipid-transfer protein 1 precursor (LTP 1) (PAPI).
ATCCATCCATCATCCATCTCATCATCAGCAACCAATTGCGACCGATCGATCGATCGATCC
...
>Os12t0114900-01 Non-protein coding transcript.
ACAAAGCGTGAGACGACAGCTGCATCTCGGCCGGCGCTGCATGTCATCCTCCGGTAGCGA
...
>Os12t0115300-01 Plant lipid transfer protein and hydrophobic protein, helical domain containing protein.
ATCCATCGCCGGGACGAGGAGACGCTGACTGAAACGAAGCTTACCTAGCTACTCGATCGA
...
>Os12t0188700-02 Similar to Thioredoxin (TRX).
ACAAAAAATTTCTCCTCCTTTCTCTTCTCTGTGTCACTGCGTGTGAGGAGCGTAG
...
>Os12t0189400-01 Similar to Photosystem I reaction centre subunit N, chloroplast precursor (PSI- N).
GTCGTCGTCTCCAACCAACCAAAATTCTTGACACGCGCAGCTCGAGCCATGGCCGGAGTG
...
```

4 解析ソフトの準備

[Bowtie ダウンロード](#) (リードをレファレンスにアライメントする)

[Samtools ダウンロード](#) (アライメント結果の整理)

[Artemis ダウンロード](#) (アライメント結果の表示用)

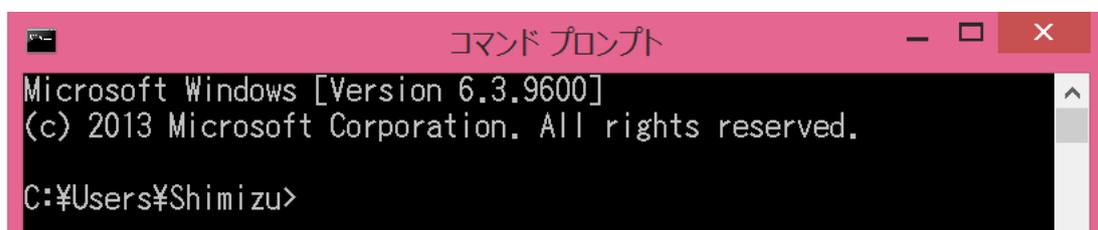
Bowtie と Samtools は圧縮ファイルを解凍します。

5 操作準備

Bowtie と Samtools は、コマンドプロンプトで操作します。その作業の便宜上、解凍してできたフォルダはそれぞれ bowtie および samtools と名前を変更し、C ドライブ直下に作成した RNAseq フォルダに移動しておきます。また、テストデータおよび参照配列データも RNAseq フォルダに移動しておきます。

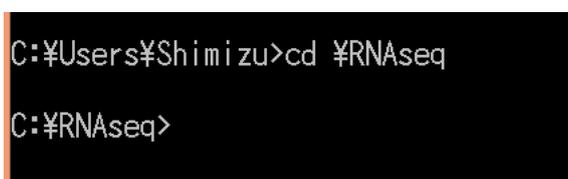
コマンドプロンプトは、コマンド入力により PC を操作するツールです。スタートメニューなどから起動すると、下図のような黒背景に白文字のウィンドウが開きます。普段の Windows の操作

と異なり、コマンドプロンプトでは全ての操作をコマンド入力で行います。>の右側の点滅しているアンダーバーは、カーソルといってキーボードから入力した文字情報はこの位置に出力されます。



コマンドプロンプトでファイルを指定するとき、その位置情報 (Path という) が重要です。上の図で、>の左側の文字列 C:¥Users¥Shimizu は現在のパスを表しており、C:はドライブ名を、¥はルートと呼びパスの区切りを表しており、¥Users¥Shimizu は Users フォルダの Shimizu フォルダにいることを表しています。コマンドプロンプトではフォルダをディレクトリと呼びます。

ファイルは C ドライブの RNAseq フォルダに置いているので、**cd ¥RNAseq** と入力して Enter キーを押して RNAseq フォルダに移動します。**cd** は change directory の意味です。



カーソルの左側の現在位置の表示が RNAseq に変わります。この状態で **dir** と入力して Enter キーを押してみると、このフォルダの中にあるファイルおよびフォルダが表示されます。



予め作成しておいた 2 個のファイルとおよび 2 個のフォルダが確認できます (ちなみに、. はカレントディレクトリを、.. はその 1 つ上の階層にある親ディレクトリを表します)。

コマンドプロンプトの操作は、Unix の操作に似ています(`dir` コマンドは Unix では `ls` コマンドになるなど、違いもある)。

6 解析手順

6.1 Bowtie によるマッピング

```
C:¥RNAseq>cd .¥bowtie
```

と入力し、カレントディレクトリ (ここでは C:¥RNAseq) の下にある bowtie フォルダに移動しておきます。次に、

```
>bowtie-build ¥RNAseq¥rap_sel.fasta rap_sel
```

と入力し、参照配列を bowtie で使用できるようにフォーマットします。すると 6 つのファイル(～.1.ebwt, ～.2.ebwt, ～.3.ebwt, ～.4.ebwt, ～.rev1.ebwt, ～.rev2.ebwt)、が作成できていれば OK です。

例データを参照配列にアライメントするには、次のようにします。

```
>bowtie -S rap_sel ¥RNAseq¥RiceRNAseq.fastq RiceEST.sam
```

```
C:¥RNAseq¥bowtie>bowtie -S rap_sel ¥RNAseq¥RiceRNAseq.fastq RiceEST.sam
# reads processed: 54466
# reads with at least one reported alignment: 9538 (17.51%)
# reads that failed to align: 44928 (82.49%)
Reported 9538 alignments to 1 output stream(s)
```

6.2 Samtools によるファイル変換

アライメントの結果を閲覧するのに Samtools を使います。sam 形式のファイルはヒトが読めるようにまとめられていて便利ですが、ファイルサイズを圧縮した方が扱いやすいので bam 形式に変換しておきます。

一つ上の階層にある samtools フォルダに移動するため、

```
C:¥RNAseq¥bowtie>cd ..¥samtools
```

と入力します。その後、

```
>samtools faidx ¥RNAseq¥rap_sel.fasta
```

```
>samtools import ¥RNAseq¥rap_sel.fasta.fai ¥RNAseq¥bowtie¥RiceEST.sam RiceEST.bam
```

```
>samtools sort RiceEST.bam RiceEST_sorted
```

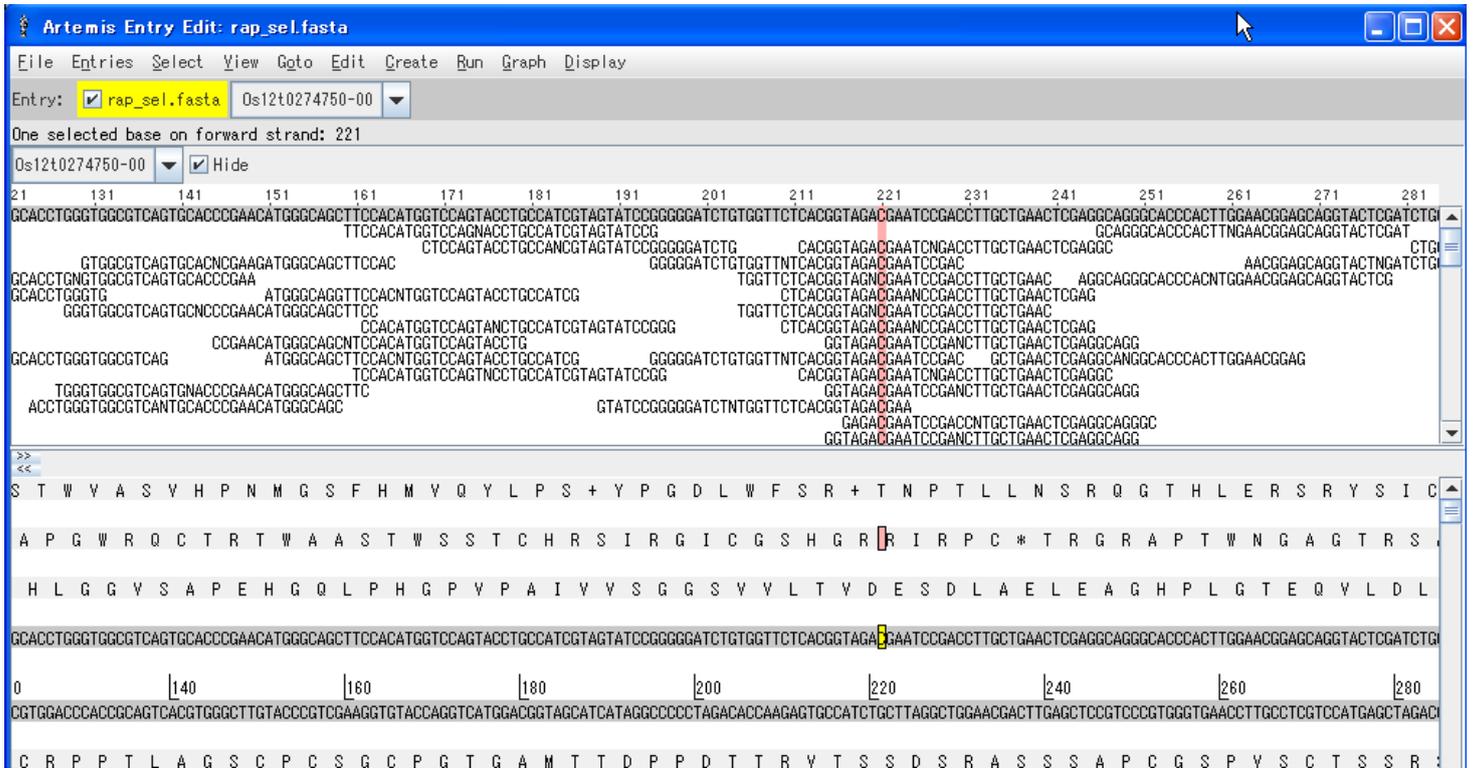
```
>samtools index RiceEST_sorted.bam
```

と入力し、出力される RiceEST_sorted.bam が閲覧可能なアライメント結果になります。

6.3 ソート済み bam ファイルの閲覧

ソート済みの bam ファイルを閲覧できるソフトウェアは複数ありますが、ここでは手軽な JAVA ソフト Artemis を使ってみます。

artemis.jar をクリックしソフトウェアを起動します。File メニューから [Open...] を選び、レファレンス配列 (FASTA 形式) を指定します。エントリー編集画面に切り替わりレファレンス配列が表示されたら、File メニューから [Read BAM/VCF...] を選び、ソート済みの BAM ファイルを指定します。



6.4 sam ファイルの中身

Bowtie が出力する sam ファイルを Excel で開いてみると、下の様になります。

まず、A 列の 2 行目から参照配列のリストが表示され (@SQ で始まる行、例データの場合 344 個ある)、

