

## 実践演習1 HISAT2/StringTie を用いた RNA-seq 解析

ここでは Nature Protocol に掲載された、HISAT2/StringTie/Balgonn を使った RNA-seq 解析手順(Pertea, M. et al. Nat. Protoc. 11, 1650, 2016)に基づき、Bias5 上でこの解析を行ってみよう。論文は、以下のサイトから取得できるので、手順の詳細については論文を参照すること。  
<https://pubmed.ncbi.nlm.nih.gov/27560171/>

1. コース用ディレクトリの下に `ex1` に移動し、その中のデータファイル `chrX_data.tar.gz` を展開する (以下、`$`で始まる行は実行するコマンドを表す)。

```
$ cd ~/bias2111/ex1
$ tar xvfz chrX_data.tar.gz
```

展開してできたディレクトリ `chrX_data` の中の以下のファイルを使用する：

`chrX_data/samples`： 入力となるリード配列 (FASTQ ファイル)

`chrX_data/indexes`： リファレンスとして用いるゲノムのインデックスファイル

`chrX_data/genes/chrX.gtf`： 遺伝子アノテーション

2. `chrX_data/samples` の下には `paired end` で 12 サンプルの配列データが入っている。このうち、とりあえず 1 サンプル分の解析を行ってみよう。以下はサンプル `ERR188044` について、HISAT2 でマッピングし、StringTie でアセンブルするという一連の解析を行う `qsub` スクリプトである。赤字で示した部分が `qsub` に特化した記述である。

`qsub_hisat2.sh`

```
#!/bin/sh
#PBS -l ncpus=4
#PBS -N hisat2
#PBS -q small

cd $PBS_O_WORKDIR
fname=ERR188044_chrX
input_dir=chrX_data/samples
seq_dir=chrX_data/indexes
gene_dir=chrX_data/genes
output_dir=output

mkdir -p $output_dir
hisat2 -p $NCPUS --dta -x $seq_dir/chrX_tran ¥
  -1 $input_dir/${fname}_1.fastq.gz -2 $input_dir/${fname}_2.fastq.gz ¥
  -S $output_dir/${fname}.sam
samtools sort -@ $NCPUS -o $output_dir/${fname}.bam $output_dir/${fname}.sam
stringtie -p $NCPUS -G $gene_dir/chrX.gtf -o $output_dir/${fname}.gtf ¥
  $output_dir/${fname}.bam
```

スクリプト先頭で PBS のオプションとして `-l ncpus=4` が指定されているので、スクリプト

内の変数`$NCPUS`に4が代入されることに注意しよう。3箇所ある`$NCPUS`は、それぞれのコマンドで使用するCPU（スレッド）数を指定しており、ここでは各コマンドが一貫して4CPUを使って動作することになる。

以下のコマンドでジョブをサブミットする

```
$ qsub qsub_hisat2.sh
```

ジョブの状態を確認する。`USERNAME`は自分のユーザ名を入れること。

```
$ qstat -u USERNAME
```

Status（後ろから2カラム目）がR（実行中）になっていることを確認しよう。ここがQの場合は実行待ち状態なので、実行開始されるまでしばらく待つ必要がある。

`qstat`で何も表示されなくなるとジョブは終了している。`ls -lt`で新たに出力されたファイルを確認しよう。正しく動作していれば、`output`ディレクトリが作成され、結果がこの中に格納されるはずである。この他に、コマンド実行時の標準出力、標準エラー出力がそれぞれ`hisat2.oXXXXXX`、`hisat2.eXXXXXX`（XXXXXXにはジョブIDが入る）というファイルに格納される。中身を見て、異常終了している様子などがいないかを確認しよう。

また、以下のコマンドを実行して、実行ログとジョブが使用したリソースを確認しよう。

```
$ tracejob XXXXXX （XXXXXXは上記のファイル名と同じジョブIDを入れる）
```

使用リソースは終了状態を示す最後の行に `resource_used.resource_type=###`の形式で表示されており、`resource_type`として`cpuct`はCPU時間、`walltime`は実時間、`ncpus`はCPU数、`mem`は使用メモリ量を表している。

3. 前項で行った解析を、`fname`で指定したサンプル名を変えつつ、12サンプル分繰り返して行う。これはfor文を使ったスクリプトを書いて実行しても良いが、ここではアレイジョブを使って実行してみよう。アレイジョブでは、スクリプト中に`$PBS_ARRAY_INDEX`という変数を埋め込み、これに指定した範囲の値を順次代入した複数のジョブを生成し、並列に実行する。アレイジョブを使うには、入力ファイル、出力ファイルともに、`file.1`、`file.2`、...のように同じ名前に番号を付加したファイル名として扱い、これらをスクリプト内で`file.$PBS_ARRAY_INDEX`のように指定するのが基本である。

今回の入力ファイルは残念ながらそのような名前になっていない。この場合、シンボリックリンクで各入力ファイルに上記形式の別名をつけるのがわかりやすい方法であるが、ここではややトリッキーだが、より簡潔に書けるbashの配列変数を使ったやり方でやってみよう。配列変数は複数の要素を持ち、各要素を添字で取り出せるようなデータ構造で、bashでは以下のようにして定義する。

```
$ Array=( fileA fileB fileC )
```

添字は 0 から始まるので、例えば最初の要素を取り出して表示するには

```
$ echo ${Array[0]}
```

とする（中括弧が必要なので注意）。

この配列変数を使って、前項のスクリプトを、アレイジョブを使うように修正したのが以下のスクリプトである。修正箇所は青字で表示されている。

qsub\_hisat2\_array.sh

```
#!/bin/sh
#PBS -l ncpus=4
#PBS -N hisat2
#PBS -q small
#PBS -J 0-11

cd $PBS_O_WORKDIR
input_dir=chrX_data/samples
seq_dir=chrX_data/indexes
gene_dir=chrX_data/genes
output_dir=output

input_files=(`ls $input_dir | grep _1.fastq | sed 's/_1.fastq.gz/'`)
fname=${input_files[$PBS_ARRAY_INDEX]}

mkdir -p $output_dir
hisat2 -p $NCPUS --dta -x $seq_dir/chrX_tran ¥
  -1 $input_dir/${fname}_1.fastq.gz -2 $input_dir/${fname}_2.fastq.gz ¥
  -S $output_dir/${fname}.sam
samtools sort -@ $NCPUS -o $output_dir/${fname}.bam $output_dir/${fname}.sam
stringtie -p $NCPUS -G $gene_dir/chrX.gtf -o $output_dir/${fname}.gtf ¥
  $output_dir/${fname}.bam
```

`input_files` が 12 個のサンプルに対応する入力ファイル名を格納した配列変数である。定義がややトリッキーだが、どんな値が入るかはバッククオートで囲まれている以下のコマンドを実行してみれば良い（バッククオートはその中のコマンドを実行してその結果をそのまま貼り付けることを指示している）。

```
$ ls chrX_data/samples | grep _1.fastq | sed 's/_1.fastq.gz/'
```

これで、配列変数には `ERR188044_chrX`、`ERR188104_chrX` 等の 12 個のサンプルに対応する名前が入ることが確認できる。次の行で、この配列から `$PBS_ARRAY_INDEX` で指定した要素の一つを取り出して `fname` に代入している。配列変数の添字は 0 から始まるので、`$PBS_ARRAY_INDEX` のとる値（`-J` オプション）を 0 から 11 の間としている点に注意。

ジョブをサブミットする

```
$ qsub qsub_hisat2.sh
```

`qstat` で状態を確認しよう。デフォルトではアレイジョブは 1 つのジョブとして表示され、

Status は B となっている。個々のジョブ (サブジョブ) の状態を確認したい場合は、`qstat -t` オプションを用いる。

```
$ qstat -u USERNAME -t
```

終了したら、`output` ディレクトリに 12 個のサンプルに対する解析結果が格納されていることを確認しよう。また、標準出力、標準エラー出力を格納したファイルがサブジョブごとに作成されていることも確認しよう。

4. ここまでで、各サンプルのリードをリファレンス配列にマッピングしてアセンブルし、サンプルごとに転写配列/遺伝子の位置を記録した GTF ファイルとして作成するところまで終わっている。このあと、サンプルごとのアセンブル結果をマージして一つのアノテーションファイルとして作成し、これを用いて各サンプルのマッピング結果から転写配列/遺伝子ごとのカウントデータを作成するステップに進むが、これを実行するスクリプトを `postproc.sh` として作成している。これを `qsub` で実行せよ。

(以下オプション) このスクリプトの後半では、サンプルごとにカウントデータを作成する処理を、`for` ループを使って行っている。

```
for bamfile in $output_dir/ERR*.bam; do
    fname=`basename $bamfile .bam`
    sample_name=`echo $fname | sed 's/_chrX//'`
    stringtie -e -B -p ${NCPUS} -G stringtie_merged.gtf ¥
    -o ${ballgown_out}/${sample_name}/${fname}.gtf $bamfile
done
```

これは、一つのジョブ内で逐次的に実行されるために効率が悪い。この部分の処理を、アレイジョブを使って行うスクリプトを作成してみよ。

(解答例は `ans` ディレクトリの下にある)

5. (オプション: R の実行) ここまでで元論文のステップ 6 までが終了し、`ballgown` ディレクトリ内に、サンプルごとに作成した転写配列/遺伝子のカウントデータが GTF ファイル、および `Ballgown` 用のテーブル形式のデータとして格納されている。これ以降を続けて行う場合、R を使った解析になる。ローカルマシン上の R を使う場合は、`ballgown` 以下のファイルを `scp` でダウンロードしてローカル環境で解析を行うこと。`Bias5` 上で R を実行することもできるが、その場合は、以下の点に注意すること。1) 足りないパッケージが自分でインストールすることができる。その際は、各自のホームディレクトリ上にインストールされる。2) グラフィックス表示を行うには `X windows` 環境がローカルに設定されていることが必要になる。その際はログイン時に `ssh -Y` オプションを指定すること。

6. (オプション: データベースの参照) 今回は、リファレンス配列とその検索用インデックスがあらかじめ用意されていたが、実際に解析する際はこれらを自分で用意する必要がある。その際は、`bias5` 上に置かれた以下のデータベースファイルを参照することができる。

データベース	パス
Illumina iGenomes	<code>/bio/db/igenomes</code>
Ensembl	<code>/bio/ftp/ensemble</code>
NCBI Genomes	<code>/bio/ftp/genomes/refseq</code> (または <code>genbank</code> )

これらを用いてリファレンス配列を用意し、`hisat2-build` でインデックスを作成するスクリプトとして `preproc.sh` を用意したので、興味がある人は内容を確認した上で実行してみよう。このスクリプトは、デフォルトでは `iGenomes` からデータを取得するが、コメントを付け直すことで `Ensembl` や `NCBI Genomes` からデータを取得するように変更できる。

## 実践演習2 アレイジョブを使った BLAST の実行

PBS 利用法の講義で、クエリ配列を分割してアレイジョブとして実行する例が出てきた。ここではそれを実際にやってみよう。コース用ディレクトリの下に **ex2** に移動しよう。クエリ配列は、講義で使ったのと同じ出芽酵母ゲノムの全遺伝子翻訳配列 `sce_prot.fasta` である。

配列ファイルを分割するにはいろいろな方法があるが、ここでは配列ファイル操作ユーティリティプログラム SeqKit の `split` サブコマンドを使って行ってみよう。

```
$ seqkit split sce_prot.fasta -p 50
```

このコマンドでは、配列ファイル中の配列を `-p` オプションで指定した数のファイルに分割する。結果は `sce_prot.fasta.split` というディレクトリに格納される。この中身を確認しよう。`sce_prot.part_nnn.fasta` という名前 (`nnn` は 001 から 050 の数値) の 50 個のファイルに分割されている。

以下のスクリプトは、これらのファイルを入力としてアレイジョブとして `blast` を実行する。

qsub\_blast.sh

```
#!/bin/sh
#PBS -l ncpus=8
#PBS -l mem=12gb
#PBS -N blast
#PBS -q small
#PBS -J 1-50

cd ${PBS_O_WORKDIR}

seqn=`printf "%03d" $PBS_ARRAY_INDEX`
query=sce_prot.fasta.split/sce_prot.part_${seqn}.fasta

db=swissprot
outdir=sce_prot.${db}.blast_out
output=$outdir/sce_prot.part_${seqn}.tab
mkdir -p $outdir

blastp -num_threads $NCPUS -db $db -outfmt 6 -query $query ¥
-evalue 0.001 -out $output
```

ただし、青字はアレイジョブに特有の記述、赤字はその他の PBS に特有の記述である。

オプション `-J` を指定しているため、変数 `$PBS_ARRAY_INDEX` に `-J` で指定した `1-50` の値を順に代入した 50 個のサブジョブからなるアレイジョブとして実行される。ただし、ここでは入力ファイル名が 0 から始まる 3 桁の数値になっているため、この形式に変換するために `printf` コマンドを用いている。これは C 言語でデータをフォーマット出力する関数 `printf` と同じ仕様のコマンドで、`"%03d"` で 0 から始まる 3 桁の整数値への変換を指示している。例え

ば以下のコマンドは（改行なしで）010 を出力する。

```
$ printf "%03d" 10
```

また、（サブ）ジョブあたりで確保するリソースとして CPU 数 8 個とメモリ 12GB を指定している。BLAST は使えるメモリがある限りメモリ上にデータベースを展開して実行されるので、大きなデータベースを使用する際はメモリを最大限確保する必要がある。このため blast キューが用意されているが、今回用いるデータベース swissprot はそこまで大きくないので、通常の small キューでも十分である。

これを qsub で実行しよう。

```
$ qsub qsub_blast.sh
```

qstat -t でサブジョブごとの実行状況を確認する。

```
$ qstat -u USERNAME -t
```

このとき、すべてのジョブが実行状態 R にはなっておらず、一部は待ち状態 Q になっているはずである。なぜか。いま、50 個のジョブがキューに入っており、ジョブあたり 8 個の CPU を確保して実行している。そこで、これらのジョブが使用する CPU 数の合計は 400 個になっており、これは small キューにおいて、ユーザあたりで最大同時に実行できる CPU 数 300 を超えている（Bias におけるキュー構成のテーブルを参照のこと）。この場合、同時に実行されるジョブ数は  $300/8 = 37.5$  を超えることができず、残りは待ち状態となる。

一般に独立したジョブを多数流す場合、ジョブあたりの CPU 数を増やすよりは、ジョブ数を増やした方が、同じ CPU 数を使っても効率が良いことが多い。したがって、このようにジョブあたりの CPU 数を増やして待ち状態が生じると、一般に実行効率は低下する。今回の場合は、ジョブあたりの CPU 数を 6 個以下にすればこのような制約による待ち状態は生じない。

上記の実行が終わったら、ファイルを編集して CPU 数を 6 に変更して (ncpus=6) 再実行してみよう。ただし、他のユーザも含めたキューの混み具合によっても同時に実行されるジョブ数は変わってくるので注意すること。