# RNA sequence data analysis
## (Part 1: using pathogen portal's RNAseq pipeline)
## Exercise 3

The goal of this exercise is to retrieve an RNA-seq dataset in FASTQ format and run it through an RNA-sequence analysis pipeline.

**Step I:** Create a login account at Pathogen Portal:

1. Go to http://pathogenportal.org
2. Click on RNA Rocket.
3. Click on Create account and fill in the required information.

**Step II:** Getting data into your launch pad.

The following exercise is based on data generated from the recent study:
Grisdale *et al.* Transcriptome analysis of the parasite *Encephalitozoon cuniculi*: an in-depth examination of pre-mRNA splicing in a reduced eukaryote. BMC Genomics 2013, 14:207

http://www.biomedcentral.com/1471-2164/14/207

In the paper the authors indicate that the data has been deposited to the sequence read archive (SRA) and a study accession number is provided: SRP017112. You can access this record here:

http://www.ncbi.nlm.nih.gov/sra/SRP017112

The required input format is something called a FASTQ file, which is similar to a FASTA file. These are simple text files that include sequence and additional information about the sequence (ie. name, quality scores, sequencing machine ID, lane number etc.).

FASTA

Definition line

>SEQUENCE_1
MTEITAAMVKELRESTGAGMMDCKNALSETNGDFDK
AVQLLREKGLGKAAKKADRLAAEGLVSVKVSDDFTIAA
MRPSYLSYEDLDMTFVENEYKALVAELEKENEERRRL
KDPNKPEHKIPQFASRKQLSDAILKEAEEKIKEELKAQ
GKPEKIWDNIIPGKMNSFIADNSQLDSKLTLMGQFYVM
DDKKTVEQVIAEKEKEFGGKIKIVEFICFEVGEGLEKKT
EDFAAEVAAQL

Sequence

FASTQ

End of Sequence

Definition line

@SRR016080.2 20AKUAAXX:7:1:123:268
TGTAGCATAATGCCGTTTTCTTTGTTTCCATTCATC
+
II&I&4IICIIIIIIII.III3:III3#6IIII1I)
@SRR016080.3 20AKUAAXX:7:1:112:638
TATAGATCTTGGTAACACCCGTTGTATTATTCGCAA
+
IIIIIIIIIIIIIIIIIIIIIIII-IIIII%%IIII
@SRR016080.4 20AKUAAXX:7:1:102:360
TTGCCAGTACAACACCGTTTTGCATCGTTTTTTTTA
+
IIIIII$IIIIIIII'IIIIIIIIIIII@IIIID35

Sequence

Encoded Quality Score

- FASTQ files are large and as a result not all sequencing repositories will store this format. However, tools are available to convert, for example, NCBI's .SRA format to FASTQ. The file that we will be using for this exercise originated from the DNA Data Bank of Japan (DDBJ), which is a mirror of NCBI and EBI.

Here is the record at DDBJ:

http://trace.ddbj.nig.ac.jp/DRASearch/study?acc=SRP017112

The FastQ files for each time point are available here:

ftp://ftp.ddbj.nig.ac.jp/ddbj_database/dra/fastq/SRA061/SRA061150/

The 24hr time data are in the folder called: SRX247417
The 48hr time data are in the folder called: SRX229331
The 72hr time data are in the folder called: SRX247418

We will be uploading data directly from the DDBJ FTP site. Each samples is paired end (ie. two files per sample). Also, they indicate that two runs were done for each sample. We are only going to worry about one of the runs for each time point. For the next part of the this exercise feel free to navigate in the FTP site to the desired time point folder or simply use the links provided below:

**Group 1 (24hr time point):**

Upstream:

ftp://ftp.ddbj.nig.ac.jp/ddbj_database/dra/fastq/SRA061/SRA061150/SRX247417/SRR769604_1.fastq.bz2

Downstream:

ftp://ftp.ddbj.nig.ac.jp/ddbj_database/dra/fastq/SRA061/SRA061150/SRX247417/SRR769604_2.fastq.bz2

**Group 2 (48hr time point):**

Upstream:

[ftp://ftp.ddbj.nig.ac.jp/ddbj_database/dra/fastq/SRA061/SRA061150/SRX229331/SRR769606_1.fastq.bz2](ftp://ftp.ddbj.nig.ac.jp/ddbj_database/dra/fastq/SRA061/SRA061150/SRX229331/SRR769606_1.fastq.bz2)

Downstream:

[ftp://ftp.ddbj.nig.ac.jp/ddbj_database/dra/fastq/SRA061/SRA061150/SRX229331/SRR769606_2.fastq.bz2](ftp://ftp.ddbj.nig.ac.jp/ddbj_database/dra/fastq/SRA061/SRA061150/SRX229331/SRR769606_2.fastq.bz2)

**Group 2 (72hr time point):**

Upstream:

[ftp://ftp.ddbj.nig.ac.jp/ddbj_database/dra/fastq/SRA061/SRA061150/SRX247418/SRR769608_1.fastq.bz2](ftp://ftp.ddbj.nig.ac.jp/ddbj_database/dra/fastq/SRA061/SRA061150/SRX247418/SRR769608_1.fastq.bz2)

Downstream:
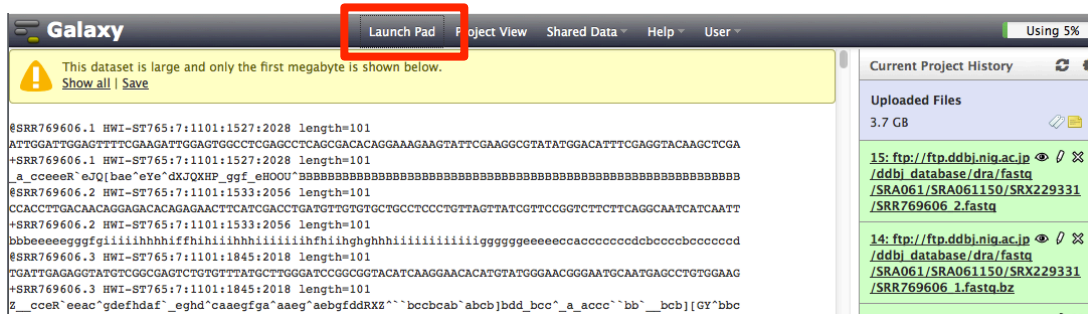
[ftp://ftp.ddbj.nig.ac.jp/ddbj_database/dra/fastq/SRA061/SRA061150/SRX247418/SRR769608_2.fastq.bz2](ftp://ftp.ddbj.nig.ac.jp/ddbj_database/dra/fastq/SRA061/SRA061150/SRX247418/SRR769608_2.fastq.bz2)

Here are the steps you take to start uploading data into your Launchpad:

1. Click on the "Upload Files" link

2.  On the next page, copy and paste both files for your time point in the "URL/Text" window then click on the "Execute" button.



You should now see a window that looks like this:



To view the progress of your upload, click on "Project View" (red square in image above).



In progress tasks will show up in yellow

Completed tasks will show up in green

You can inspect the contents of completed tasks (like uploaded files) by clicking on the eye icon next to the name of the file (arrow in above image).  Inspecting a FASTQ file should look like this:



3. Once the RNA-sequence FASTQ file has been uploaded you can start the RNA-seq pipeline.  Pathogen portal uses two algorithms for mapping (TopHat) and transcript prediction and expression value calculation (Cufflinks).  Note that there are many algorithms and methods for RNA-seq mapping and analysis each with its advantages and disadvantages.  You are encouraged to learn more about the algorithm you are using.

   o  TopHat:      http://tophat.cbcb.umd.edu/
   o  Cufflinks:   http://cufflinks.cbcb.umd.edu/index.html

- To start the pipeline click on the "Launch Pad" link (red square in above image). On the next page, scroll down to the "RNA-Seq Analysis" section and click on "Align Reads & Assemble Transcripts".

- On the next page, scroll down and choose the type of analysis (in this case we are analyzing a paired end eukaryotic sample).
- Next select the target project from the drop down menu. You should only have one or two projects one of which will contain both FASTQ files you uploaded (probably called "Uploaded Files"). Once you select the correct project you should see the two FASTQ files contained within it. Next click on continue.



- The next page allows you to configure the pipeline:

  **Step1:** Select the upstream read file (ends in _1) and click on the arrow to move it to the "Selected" window.

  **Step2:** Select the downstream read file (ends in _2) and click on the arrow to move it to the "Selected" window.

**Step3:** Configure TopHat – there are a number of options that may be modified, however, for the purposes of this exercise the default parameters may be used. The only required change is the reference genome -- select *Encephalitozoon cuniculi* EC2

**Step 3: Tophat2** (version 2.0.6)

**Is this library mate-paired?**
Paired-end

**RNA-Seq FASTQ file, forward reads**
Output dataset 'output' from step 1
Nucleotide-space: Must have Sanger-scaled quality values with ASCII offset 33

**RNA-Seq FASTQ file, reverse reads**
Output dataset 'output' from step 2
Nucleotide-space: Must have Sanger-scaled quality values with ASCII offset 33

**Mean Inner Distance between Mate Pairs**
300

**Std. Dev for Distance between Mate Pairs**
20
The standard deviation for the distribution on inner distances between mate pairs.

**Report discordant pair alignments?**
No

**Use a built in reference genome or own from your history**
Use a built-in genome
Built-in genomes were created using default options

**Select a reference genome**
Encephalitozoon cuniculi EC2
If your genome of interest is not listed, contact the Pathogen Portal team

**TopHat settings to use**
Use Defaults
You can use the default settings or set custom values for any of Tophat's parameters.

**Specify read group?**
No

**Step4:** Configure Cufflinks – once again there are a number of options to modify. For the purposes of this exercise change the following:
Maximum Intron Length (-I): 1000
Select a reference annotation: *Encephalitozoon cuniculi* EC2
Select how to use the provided annotation: Assemble Novel + annotated transcripts.

**Click on the Run Workflow button.**

**Step 4: Cufflinks Eukaryotic** (version 2.0.2)

**SAM or BAM file of aligned RNA-Seq reads**
Output dataset 'accepted_hits' from step 3

**Maximum Intron Length (–I)**  ⓘ
1000

**Minimum Isoform Fraction (–F)**  ⓘ
0.1

**Pre MRNA Fraction (–j)**  ⓘ
0.15

**Overlap Radius**  ⓘ
50

**Perform Quartile Normalization**  ⓘ
No

**Will you select a reference annotation from your history or use a built-in file from Pathogen Portal?**
Use provided annotation

**Select a reference annotation**
Encephalitozoon cuniculi EC2
If your annotation of interest is not listed, contact Pathogen Portal team.

**Select how to use the provided annotation**
Assemble ONLY transcripts matching the annotation

**Perform Bias Correction**
Yes
Bias detection and correction can significantly improve accuracy of transcript abundance estimates.

**Reference Sequence Data**
Locally cached

**Use multi-read correct**  ⓘ
No

None

**Run workflow**

After you start the workflow you should get a confirmation window that indicates all the steps that have been added to the queue. The progress of your workflow can be viewed to the right. Completed tasks are in green, running tasks are in yellow and tasks waiting in the queue are in grey.