# Mapping RNA sequence data
## (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.

This exercise will rely on data deposited in the sequence read archive (SRA).  The data is based on transcriptomic analysis of three developmental stages of *Plasmodium falciparum*: 1. Cultured asexual stages, 2. Cultured sporozoites, and 3. Salivary gland Sporozoites. Two replicates of each developmental stage were sample were generated. Additional information about this experiment may be obtained from GEO:

http://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE52867

Examining the information available in GEO and under the SRA accession numbers you will notice that this data is paired end.  So for each sample there should be two files one for each of the pairs.

Salivary gland sporozoites sample 1:	http://www.ncbi.nlm.nih.gov/sra/SRX385640
Salivary gland sporozoites sample 2:	http://www.ncbi.nlm.nih.gov/sra/SRX385641
Cultured sporozoites sample 1:	http://www.ncbi.nlm.nih.gov/sra/SRX385642
Cultured sporozoites sample 2:	http://www.ncbi.nlm.nih.gov/sra/SRX385643
Asexual stage parasites sample 1:	http://www.ncbi.nlm.nih.gov/sra/SRX385644
Asexual stage parasites sample 2:	http://www.ncbi.nlm.nih.gov/sra/SRX385645

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

>SEQUENCE_1

MTEITAAMVKELRESTGAGMMDCKNALSETNGDFDK
AVQLLREKGLGKAAKKADRLAAEGLVSVKVSDDFTIAA
MRPSYLSYEDLDMTFVENEYKALVAELEKENEERRRL
KDPNKPEHKIPQFASRKQLSDAILKEAEEKIKEELKAQ
GKPEKIWDNIIPGKMNSFIADNSQLDSKLTLMGQFYVM
DDKKTVEQVIAEKEKEFGGKIKIVEFICFEVGEGLEKKT
EDFAAEVAAQL

Sequence

## FASTQ

Definition line

End of Sequence

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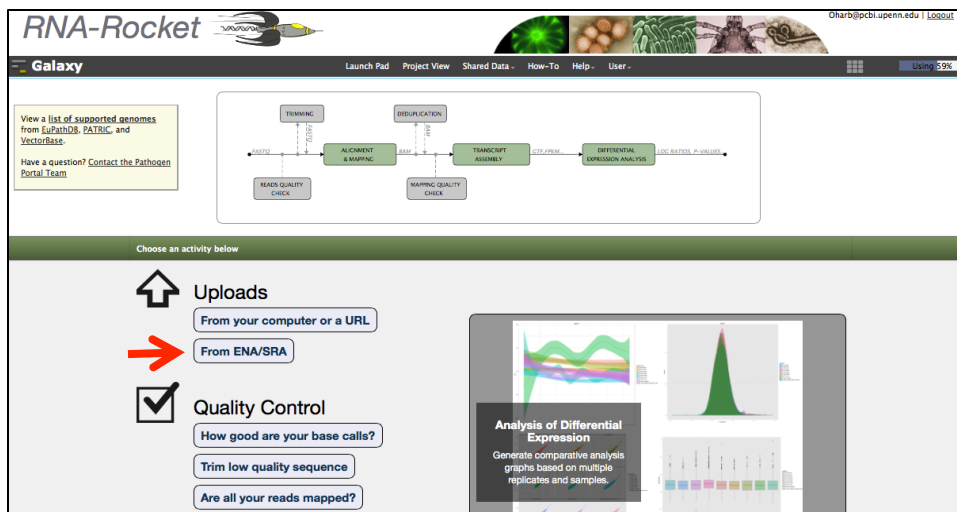IIIIII$IIIIIIII'IIIIIIIIIIII@IIIID35
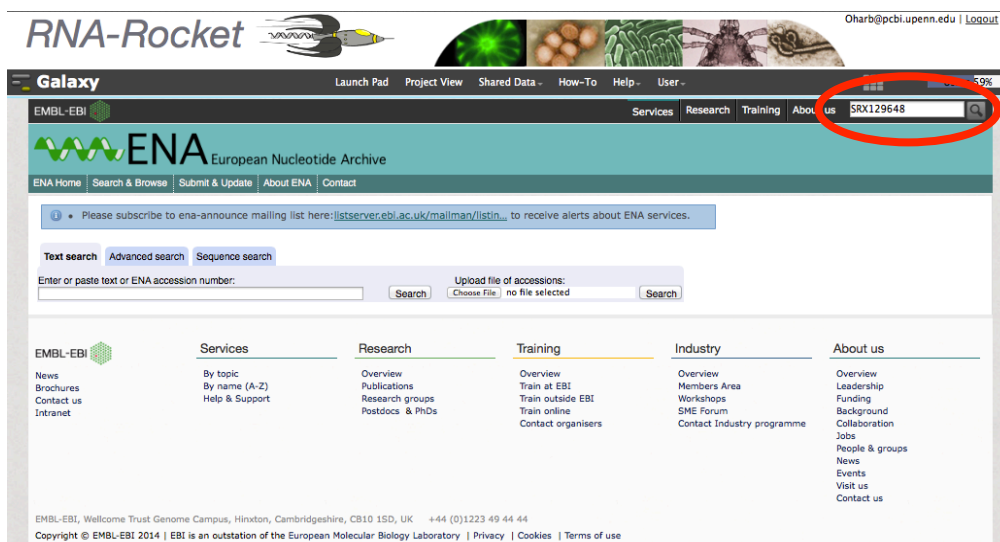
Sequence

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.
- Sequence data is housed in three repositories that are synchronized on a regular basis.

  o The sequence read archive at GenBank
  o The European Nucleotide Archive at EMBL
  o The DNA data bank of Japan

- RNArocket allows you to use SRA accession numbers and directly retrieve FASTQ files.

➢ Here are the steps you take to start uploading data into your Launchpad:
   o Note: During this exercise you will <u>NOT</u> download any data to your computer. Instead you will be providing information to enable transferring data from SRA to RNA-Rocket.

1. Click on the "Upload Files" link
2. On the next page, notice the instructions to use the global search on the ENA site. Next click on continue.



3. Cut and paste the study accession number (SRP033414) into the global search box (see red circle below). Click on the search icon.

4. Click on the Study link obtained.



5. To transfer files to RNA-Rocket, click on the File 1 or File 2 in the column called "Fastq files (galaxy)". Remember, you have to get 2 files, one for each pair. Click on the link for File 1 for the sample assigned to your group, then click on the back button on your browser and click on the link for File 2 from the same sample.

You should now see a window that looks similar to 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

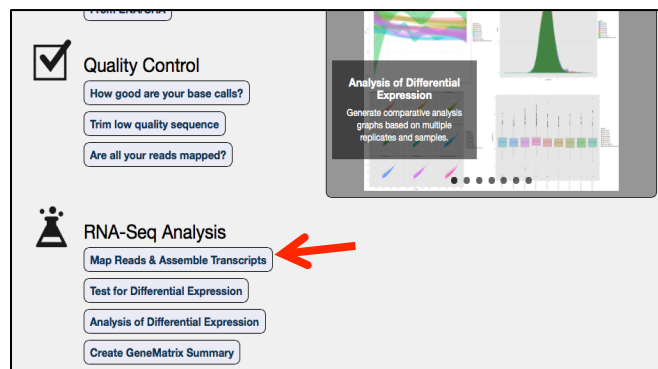You can inspect the contents of completed tasks (like uploaded files) by clicking on the eye icon next to the name of the file (arrow in above image).  Inspecting a FASTQ file should look like this:

6. 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 "Map 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 *Plasmodium falciparum* 3D7



Step 3: Tophat2 (version 2.0.10)

**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?**
Yes

**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**
Plasmodium falciparum 3D7
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): 5000
The reference annotation should be automatically selected: *Plasmodium falciparum* 3D7
Select how to use the provided annotation: Assemble Novel + annotated transcripts.

**Click on the Run Workflow button.**



Step 4: Cufflinks (version 2.0.2)

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

**Maximum Intron Length (-I)**
5000

**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**
Plasmodium falciparum 3D7
If your annotation of interest is not listed, contact Pathogen Portal team.

**Select how to use the provided annotation**
Assemble novel+annotated transcripts

**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.