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



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

1. Go to http://pathogenportal.org
2. Click on Analyze at the top of the page (A).
3. Click on "RNA-Seq Pipeline (B).
4. Click on Create account and fill in the required information (C).

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

For this exercise we will be working with a data set generated from Illumina sequencing of a *Babesia bovis* cDNA library (http://www.ncbi.nlm.nih.gov/pubmed/18987005).

You can read more about the actual sample files here:
http://www.ncbi.nlm.nih.gov/sra/SRX004534

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
+
IIIIIIIIIIIIIIIIIIIIIIIII-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 NCBI:

http://www.ncbi.nlm.nih.gov/sra/SRX004534

Here is the record at DDBJ:

http://trace.ddbj.nig.ac.jp/DRASearch/run?acc=SRR016080

- To save download and upload time the file is available as a shared project in PathogenPortal. Go to the following link, then click on "Import History".

http://rnaseq.pathogenportal.org/u/omar/h/babesiabovis



- Once the RNA-sequence FASTQ file has been imported into your history 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.
    - TopHat:       http://tophat.cbcb.umd.edu/
    - Cufflinks:    http://cufflinks.cbcb.umd.edu/index.html

- To start the pipeline click on the "Launch Pad" link.
  - Select the file you just imported.
  - Choose the workflow – in the case were running a "Eukaryotic Single End Analysis".
  - Choose a destination project.  You can give this a name in the "New Project Named" window.
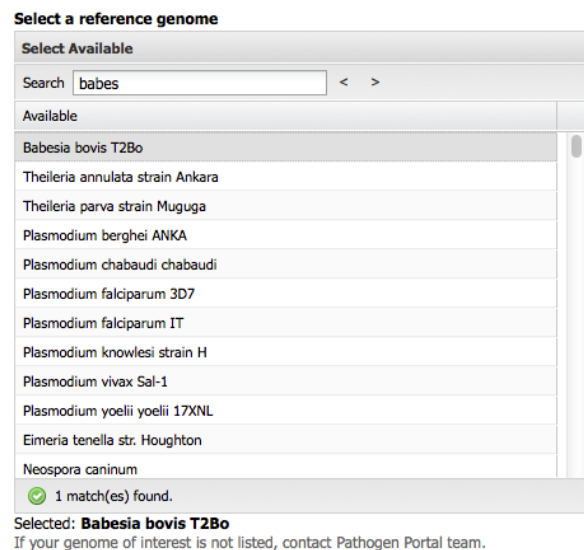  - Click on Continue.



- The next page allows you to configure the pipeline and set both TopHat and Cufflinks parameters:
  - Step1: Select input dataset.
  - Step2: Configure TopHat.
  - Step3: Configure Cufflinks.

Step 1: Select input dataset and click on the arrow to move it from the "Available" window to the "Selected" window:



Step 2: Configure TopHat:

- o Select the reference organism – in this case *Babesia bovis*. You can start typing the name of the organism in the search window. This will automatically search for the closest match and select it for you.
- o We will leave most of the TopHat paramaters at the default values.
- o Change the "Maximum intron length" field to 2000.





Step 3: Configure cufflinks:

- o Change the "Maximum intron length" field to 2000.
- o Select the reference annotation – in this case *Babesia bovis* (exactly as you did above).

Click on the Run Workflow button. 

After you start the workflow you should get a confirmation window that indicates all the steps that have been added to the queue:



Successfully ran workflow "Eukaryotic Single-End Analysis". The following datasets have been added to the queue:
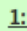
1: SRR016080.fastq

2: Tophat for Illumina on data 1: insertions

3: Tophat for Illumina on data 1: deletions

4: Tophat for Illumina on data 1: splice junctions

5: Tophat for Illumina on data 1: accepted_hits

6: Cufflinks on data 5: gene expression

7: Cufflinks on data 5: transcript expression

8: Cufflinks on data 5: assembled transcripts

9: Cufflinks on data 5: total map mass

You can check the progress of your workflow by clicking on the "Project View" link. Completed tasks are in green, running tasks are in yellow and tasks waiting in the queue are in grey: