

Interpreting RNAseq Mapping results

(Part 2: Loading data generated by the pathogen portal's RNAseq pipeline in the Genome Browser)

Exercise 7

For this exercise we will be using:

<http://pathogenportal.org>

<http://giardiadb.org>

1. Explore the results of the RNA-sequence pipeline. What files were generated? To view contents of any of the results, click on the eye icon () next to the file name.

!!! important note - do not click on the icon next to the file called "Tophat2 on data 1 and data 3: accepted_hits" - this file is huge and will not display but rather will download the contents to your computer.

- a. TopHat in RNArocket generates five files:
Align_summary: this includes a summary of how the alignment went (ie. the number of reads that were aligned).
- b. *Insertions*: reported insertions.
- c. *Deletions*: reported deletions.
- d. *Splice junctions*: reported junctions. Each junction consists of two connected BED blocks, where each block is as long as the maximal overhang of any read spanning the junction. The score is the number of alignments spanning the junction.
- e. *Accepted hits*: BAM file (binary alignment map).
Note that many alignment programs will generate a file called a SAM file (sequence alignment map) which is a table including text of the alignment and mapping. However, for viewing results in a sequence browser like GBrowse, the file needs to be converted into the binary formatted (BAM) - you do not have to worry about this for this exercise.

14: Tophat2 on data 2 and data 1: accepted hits (Genome Coverage BedGraph)   
13: Tophat2 on data 2 and data 1: accepted hits (- BigWig)   
12: Tophat2 on data 2 and data 1: accepted hits (+ BigWig)   
10: Cufflinks on data 7: assembled transcripts   
9: Cufflinks on data 7: transcript expression   
8: Cufflinks on data 7: gene expression   
7: Tophat2 on data 2 and data 1: accepted hits   
6: Tophat2 on data 2 and data 1: splice junctions   
5: Tophat2 on data 2 and data 1: deletions   
4: Tophat2 on data 2 and data 1: insertions   
3: Tophat2 on data 2 and data 1: align_summary   
2: EBI SRA: SRX129648 File: ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR445/SRR445171/SRR445171_2.fast   

Cufflinks generates three files:

gene expression, transcript expression and assembled transcripts. The gene

expression and transcript expression files for our purposes should be identical since EuPathDB genomes do not have separate genes and transcripts. These files include the FPKM values (Fragments Per Kilobase of transcript per Million mapped reads) for each gene in the genome analyzed - in this case *Giardia* assemblages.

Additional files include files of the format BigWig and BedGraph. You can read more about these file formats here:

<http://genome.ucsc.edu/goldenPath/help/bigWig.html>

In a nutshell, these are file formats created from large binary files like BAM files and makes it possible to load these data in a genome browser.

2. Load your BAM data (accepted hits) into GBrowse.

Click on your “Tophat2 on data 2 and data 1: accepted_hits” in your project history panel. This will show you information about the file including a link to display data in GiardiaDB - click on the link.

3. Load the assembled transcript data. Essentially use a similar procedure as above.

4. Wait a couple of minutes for GBrowse to load your data.

5. Once data has been loaded, you can configure the track display settings. For example, you can adjust the Y-axis scaling to a fixed axis.

The screenshot shows a list of tracks in a GBrowse interface. Track 7 is highlighted with a red box and contains the following information:

- Track title: [7: Tophat2 on data 2 and data 1: accepted_hits](#)
- Size: 1.1 GB
- Format: bam, database: gassAWB
- Log: tool progress Log: tool progress
- Timestamp: [2014-05-22 20:00:32] Beginning TopHat run (v2.0.10)
- Log details: [2014-05-22 20:00:32] Checking for Bowtie Bowtie version: 2.1.0.0 [2014-05-22 20:00:32] Checking for Samtools
- Display options: [display at EupathDB giardiadb](#) (indicated by a red arrow)
- File type: Binary bam alignments file

Other tracks visible in the list include:

- 13: Tophat2 on data 2 and data 1: accepted_hits (- BigWig)
- 12: Tophat2 on data 2 and data 1: accepted_hits (+ BigWig)
- 10: Cufflinks on data 7: assembled transcripts
- 9: Cufflinks on data 7: transcript expression
- 8: Cufflinks on data 7: gene expression
- 6: Tophat2 on data 2 and data 1: splice junctions
- 5: Tophat2 on data 2 and data 1: deletions

