RNA-seq data analysis.
Hands on workshop: RNA-seq using Galaxy.
Instructor: Wibowo Arindrarto, Hailiang (Leon) Mei, Jeroen F.J. Laros.
Leiden Genome Technology Center, Leiden University Medical Center, The Netherlands
Netherlands Bioinformatics Centre

**Introduction** In this workshop we will first show you a typical analysis done by a bioinformatician working with RNA-seq data. This involves quality control, aligning raw sequencing data to a known reference genome, doing expression analysis and visualisation using the UCSC genome browser and/or Trackster.

**Availability and examples** The tools used in these exercises are all free for download, FastQC for quality control, Sickle for data cleaning, Tophat for alignment and Cufflinks for expression analysis.

**Note on test data** Data used in this practical is test data and not full size files. This is to reduce the time needed to run each step and make this analysis possible within the time permitted. The data was retrieved from `https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE37918`, a description of the study can be found here `https://www.ncbi.nlm.nih.gov/pubmed/23580553`.

**Preparations**

- Open a browser and go to `http://ec2-72-44-50-10.compute-1.amazonaws.com/`
- Register to gain access to data libraries and workflows.
- Import the "RNA-seq Per Sample" and "RNA-seq Diff" workflow.


**Exercise 1: Single sample expression analysis** The input data is a small selection of reads that should align mostly to a small region on the human genome. After alignment, you can do expression analysis and visualisation.

Import the following files from the "GCC 2013 - RNA-seq - inputs" data library:

- `refseq_genes.gtf`
- `xtra_treat_1.fq`
- `xtra_treat_2.fq`

First look at one of the FASTQ files. Each read is represented by four lines: a header, the read itself, a "+" and the quality scores.

Do some standard QC on the FASTQ files:

- Run *FastQC* on both FASTQ files.

When looking at the output of the QC steps, you will notice a lot of warnings and errors, they arise partially from the fact that we work with a very small dataset.

*Questions*:

- Are there any other reasons for these warnings?
- What is the total number of sequences?
- What is the quality encoding?

Use *Sickle* for trimming low quality parts of the reads.

*Questions*:

- What do you see when you look at the newly generated FASTQ files?

- If you run *FastQC* again, which metrics are improved?

Align the trimmed reads to the human reference genome build `hg19` with *Tophat*.

- *Hint*: The data type should be fastqsanger

Visualise the aligned reads (BAM file) with the UCSC genome browser. Go to an area of interest. Note that splice junctions are most likely in an area of interest.

*Questions*:

- Can you find evidence for alternative splicing in region `chr22:31587843-31708266`?
  - *Hint*: Change the visualisation from "dense" to "pack".
- Can you find mismatches in the alignment (or possibly even variants)?

Run *SAMTools flagstat* on the aligned reads.

*Question*: How many reads were aligned?

- *Hint*: All of them is not the correct answer.

Make a BedGraph from the aligned reads.

- We are not interested in zero coverage regions.
- We need to take split reads into account.

*Question*: What do you see in a region of interest?

- *Hint*: Change the visualisation of the BAM- and the BedGraph track to "squish".

Inspect the insertion size metrics with *Picard* tools.

*Questions*:

- Can you explain the truncation at the left of the histogram?
- How could this be improved?

Use *Cufflinks* for transcript assembly and abundance estimation. Use the reference genes as guide for the assembly.

*Questions*:

- What is the most abundant gene?
  - *Hint*: Use the filter and sort tools.
- What is the most abundant transcript? Visualise it in the genome browser.

Extract a workflow, create a new history and apply the workflow on files:

- `refseq_genes.gtf`
- `xtra_ctrl_1.fq`
- `xtra_ctrl_2.fq`

**Exercise 2: Differential expression analysis.** Now we have analysed two samples, one treated- and one control. Now we can do differential expression analysis to figure out what the effect of the treatment was.

Import the following files from the "GCC 2013 - RNA-seq - inputs" data library:

- `ctrl.gtf`
- `treat.gtf`
- `refseq_genes.gtf`
- `ctrl.bam`
- `treat.bam`

Merge the control and treated transcript assemblies with *Cuffmerge*. Use the refseq genes as reference annotation.

Run *Cuffdiff* on the merged transcripts file, the control- and treated BAM files, use the "blind" dispersion estimation method.

For now, we are interested in the gene differential expression testing dataset. Filter this list based on the status column.

*Questions*:

- Which gene is most affected?
- Is it up- or down regulated?