# *RNA sequence data analysis via Galaxy, Part I Uploading data and starting the workflow (Group Exercise)*

The goal of this exercise is to use a Galaxy workflow to analyze RNA sequencing data. Galaxy is an open, web-based platform for data intensive biomedical research. Galaxy allows you to perform, reproduce, and share complete analyses without the use of command line scripting. EuPathDB developed its own Galaxy instance in collaboration with Globus Genomics. Many resources are available to learn how to use Galaxy. The following link has information about additional resources to help you learn how to use Galaxy:

https://wiki.galaxyproject.org/Learn#Galaxy_101

For this exercise, we will retrieve raw sequence files from a repository, assess the quality of the data, and then run the data through a workflow (or pipeline) that will align the data to a reference, calculate expression values and determine differential expression. Part 1, uploading data and starting the workflow will be performed today. The workflows will run overnight and we will view / interpret the results tomorrow in Part 2.

We will be working in groups. Each group will have 4-6 members. One person in the group will run the Galaxy controls on one computer. The other members' roles are to ensure that the correct datasets are used and that the correct workflow parameters are selected.

## Section I: Setting up your EuPathDB Galaxy account

Step 1: Access the EuPathDB Galaxy instance at the following URL:

### http://eupathdbworkshop.globusgenomics.org/

Step 2: On the next page you will be asked to define your organization. Choose EuPathDB and click Continue.

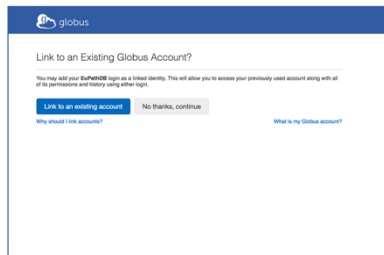Step 3: Log in to EuPathDB (if you are not logged in already).



Step 4: Next, sign up for the EuPathDB Galaxy instance.
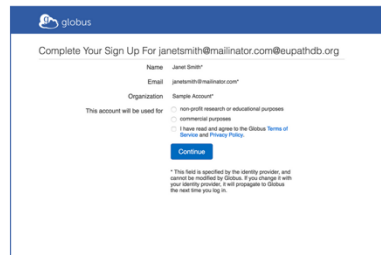
## Analyze My Experiment

The first time you visit EuPathDB Galaxy you will be asked to sign up with Globus, EuPathDB's Galaxy instance manager. This is a three-step sign-up process (screenshots below).

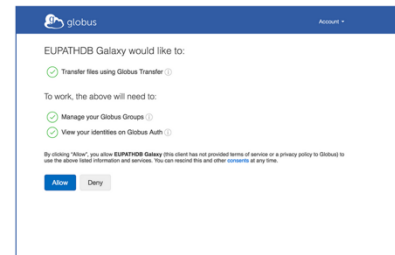Click **"Continue to Galaxy"** to sign up for EuPathDB Galaxy services.

Contact us if you experience any difficulties.



(1) If you already have a Globus account, you can link it to your EuPathDB account. **Your choice.** If you don't have a prior Globus account, choose **No Thanks.**

(2) Complete your account information and agree to Globus's Terms and Conditions. Please read, make your selections, and click **Continue.**
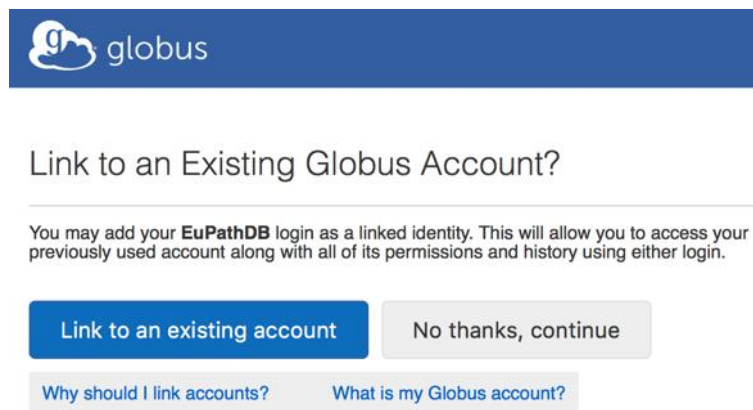
(3) Grant permission to share your Globus identity and files with us. Please click **Allow.** (We will only perform file transfers that you explicitly request, between Galaxy and other resources, including EuPathDB.)
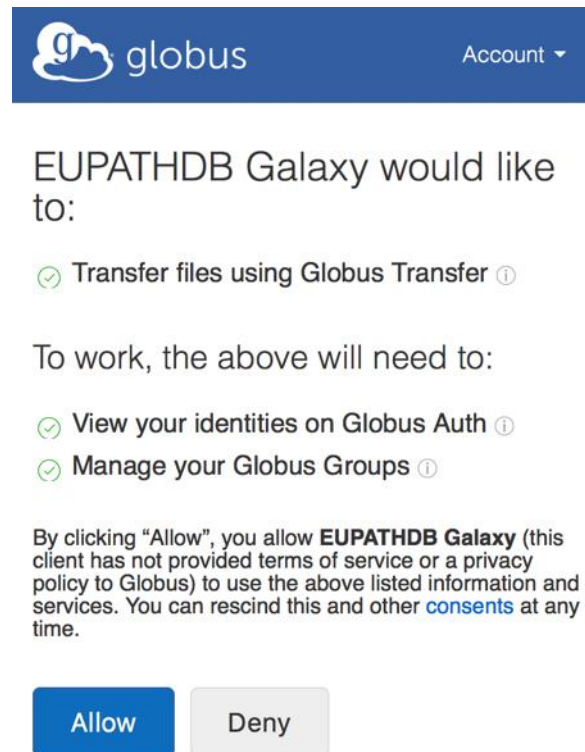
Continue to Galaxy

Step 5: Click on "Continue to Galaxy" and follow the instructions.
Step 6: Click on "No thanks, continue"

Step 7: Click on "Allow"



Step 8: Congratulations, you are in!

## *Section II: Importing data to Galaxy*

There are multiple ways to important data into your Galaxy workspace.  For this exercise, we will use the 'Download from web or upload from disk' tool and enter the direct data repository links listed below under 'Group Assignments'.  Remember one person in your group will be starting the workflow.  Although all group members can sign up for an account for later use, please only one person should start a workflow today because we don't want to overload the servers.  The samples below were all generated by paired end sequencing, hence there are two files for each sample. The files are fastq files that are compressed (that is why they end in .gz = gzip).

Step 1: Click on the "Get data" link in the left-hand menu.  This will reveal a list of options; click on "**Get Data via Globus from the EBI server**"



Step 2: In the middle section enter the sample ID and choose whether the run was single or paired end.  Click on Execute.

Note that the sample ID resulted in importing two files one for each pair. Repeat this process for each sample you want to import. *If you are working with samples from two conditions and the experiment was done in triplicate and paired end sequenced then you should end up with 12 files; six from each condition.*

Step 3: If you are working with a dataset with biological replicates it is useful to organize the different conditions of your experiment into "Collections". For example, if your experiment included RNAseq from *Plasmodium falciparum* male gametocyte stages (three biological replicates) and erythrocytic stages (three biological replicates), it is useful to organize these into two collections, one that includes all male gametocyte files and the other that includes all the erythrocytic stage files. Using collections also reduces the complexity of the Galaxy workflows. See below:

1. Click on the checkbox function "operation on multiple datastest"

2. Select samples that belong to the same condition

2. Click on "For all selected" and choose "Build a list of Datasets Pairs"

4. Usually the correct pairs are auto-selected. Double check this and give each pair a meaningful name. To change the name, click on the paired name in the center and rename it

5. Once you are done renaming the pairs, give the collection a meaningful name – for example, use the condition name. Then click on Create List

## *Group assignments:*

*Groups 1, 2 & 3* will be examining data from a study called "*Plasmodium berghei* transcriptome for female gametocytes, male gametocytes, and asexual erythrocytic stages"
https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5604118/
The data is available in the sequence repositories:
https://www.ebi.ac.uk/ena/data/view/PRJNA374918

Samples:

**Erythrocyte stages (Asexual):**
SAMN06339669
SAMN06339670
SAMN06339671

**Male gametocytes:**
SAMN06339666
SAMN06339667
SAMN06339668

**Female gametocytes:**
SAMN06339663
SAMN06339664
SAMN06339665

**Group 1:**

*Plasmodium berghei* <u>male gametocytes</u> vs. <u>erythrocytic stages</u>

**Group 2:**

*Plasmodium berghei* <u>female gametocytes</u> vs. <u>erythrocytic stages</u>

**Group 3:**

*Plasmodium berghei* <u>male gametocytes</u> vs. <u>female gametocytes</u>

***Groups 4, 5 & 6*** will be examining data from a study called "*Plasmodium falciparum* NF54 Transcriptome" which examines RNAseq from 3 stages: erythrocytic, salivary gland and cultured sporozoite stages. This study is unpublished but data is accessible in the sequence repositories: https://www.ebi.ac.uk/ena/data/view/PRJNA230379

<u>Samples:</u>

**Erythrocytic stages (Asexual):**
SAMN02428730
SAMN02428734

**Salivary gland sporozoites:**
SAMN02428726
SAMN02428729

**Cultured sporozoites:**
SAMN02428728
SAMN02428727

**Group 4:**

*Plasmodium falciparum* <u>salivary sporozoites</u> vs. <u>erythrocytic stages</u>

**Group 5:**

*Plasmodium falciparum* <u>cultured sporozoites</u> vs. <u>erythrocytic stages</u>

**Group 6:**

*Plasmodium falciparum* <u>salivary sporozoites</u> vs. <u>cultured sporozoites</u>

## Section II: Running a workflow in Galaxy

You can create your own workflows in galaxy based on your needs. The tools in the left section can all be added and configured as steps in a workflow that can be run on appropriate datasets. For this exercise we will use a preconfigured workflow that does the following main things:
1. Analyzes the reads in your files and generates FASTQC reports.
2. Trims the reads based on their quality scores and adaptor sequences (Trimmomatic).
3. Aligns the reads to a reference genome using HISAT2 and generates coverage plots.
4. Determines read counts per gene (HTSeq)
5. Determines differential expression of genes between samples (DESeq2).

- To use one of the EuPathDB preconfigured workflows, go to the Galaxy home page and select the workflow that you would like to run.  For this exercise "**EuPathDB Workflow for Illumina paired-end RNA-seq, biological replicates**" – click on this workflow to run it



- Configure your workflow – there are multiple steps in the workflow but you do not need to configure all of them.  For the purpose of this exercise you will need to configure the following:

a. Select the input dataset collections.  These are the collections of fastq files you just created.  Workflow steps 1-2 allow you to select the datasets.

b.  Some tools in the workflow require that you select the reference genome to be used. In this workflow both HISAT2 and HTSeq require this (note these tools are in the workflow twice since you have two collections). It is critical that you select the correct genome that matches the experimental organism.  So, for example, if your experiment was performed using *Plasmodium berghei*, the reference genome you select should be *Plasmodium berghei*.

**Source for the reference genome to align against**

Use a built-in genome

**Select a reference genome**

    PlasmoDB-29_Pchabaudichabaudi_Genome
    PlasmoDB-29_PcynomolgiB_Genome
    PlasmoDB-29_Pfalciparum3D7_Genome
    PlasmoDB-29_PfalciparumIT_Genome
    PlasmoDB-29_PknowlesiH_Genome
    PlasmoDB-29_PreichenowiCDC_Genome
    PlasmoDB-29_PvivaxP01_Genome
    PlasmoDB-29_PvivaxSal1_Genome
    PlasmoDB-29_Pyoeliiyoelii17XNL_Genome
    PlasmoDB-29_PyoeliiyoeliiYM_Genome
    PlasmoDB-30_PcoatneyiHackeri_Genome
    PlasmoDB-30_PfragileNilgiri_Genome
    PlasmoDB-30_PinuiSanAntonio1_Genome
    PlasmoDB-30_PmalariaeUG01_Genome
    PlasmoDB-30_PvinckeipetteriCR_Genome
    PlasmoDB-30_Pvinckeivinckeivinckei_Genome
    PlasmoDB-30_Pyoeliiyoelii17X_Genome
    PlasmoDB-32_PbergheiANKA_Genome
    PlasmoDB-32_Pgallinaceum8A_Genome
    PlasmoDB-32_PovalecurtisiGH01_Genome

**Paired alignment parameters**

c.  Another very important parameter to check in the htseq-count step is the Feature type.  The default is usually set to exon.  Make sure you chance this to gene.  To change this to gene, click on the edit icon, the type the word "gene".  This is case sensitive so be careful about this.

🔧 **htseq-count – Count aligned reads in a BAM file that overlap features in a GFF file (Galaxy Version HTSEQ: default; SAMTOOLS: 1.2; PICARD: 1.134)**

**Aligned SAM/BAM File**

Output dataset 'output_alignments' from step 7

✎ **Is this library mate-paired?**

paired-end

**Will you select an annotation file from your history or use a built-in gff3 file?**

Use a built-in annotation

**Select a genome annotation**

    PlasmoDB-32_PbergheiANKA_Genome                                          ▼

✎ **Mode**

Union

✎ **Stranded**

Yes

✎ **Minimum alignment quality**

0

⟳ **Feature type**

gene

Feature type (3rd column in GFF file) to be used. All features of other types are ignored. The default, suitable for RNA-Seq and Ensembl GTF files, is exon.

✎ **ID Attribute**

ID

d.  Once you are sure everything is configured correctly, click on "Run Workflow" at the top.

The steps will start running in the history section on the right. Grey means they are waiting to start. Yellow means they are running. Green means they have completed. Red means there was an error in the step.

**Appendix:**

FASTQ file are text files (similar to FASTA) that include sequence quality information and details in addition to the sequence (ie. name, quality scores, sequencing machine ID, lane number etc.). 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.

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