# RNA sequence data analysis via Galaxy, Part I Uploading data and starting the workflow (Group Exercise)
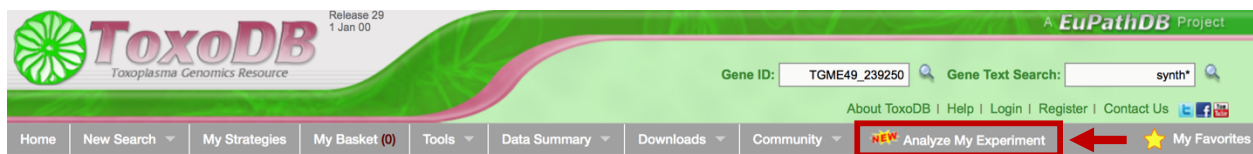
The goal of this exercise is to use a Galaxy workflow to analyze RNA sequencing data. The datasets we will us are all feely available through the sequence repositories (SRA, ENA, DDBJ).

Galaxy is an open, web-based platform for data intensive biomedical research. EuPathDB is developing its own Galaxy instance that will become available for its users. Galaxy allows you to perform, reproduce, and share complete analyses.

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 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' role is to help ensure that the correct datasets are being used and that the correct workflow parameters are being selected. You can access the EuPathDB Galaxy instance through the "Analyze My Experiment" link in the gray menu bar in any EuPathDB site:



## Section I: Setting up your EuPathDB Galaxy account

Step 1: Click on the "Analyze My Experiment" link in the gray menu bar.
Step 2: On the next page you will see a description of this service. In order to start using the EuPathDB Galaxy instance you will have to follow the registration steps. Start by clicking on the "Continue with Galaxy Sign-up" button.

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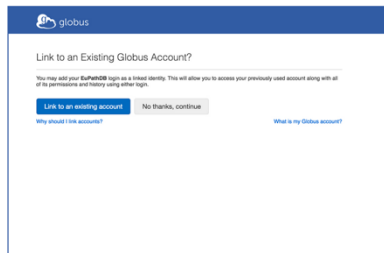


Step 4: The next window describes the process of signing up for the EuPathDB Galaxy instance.
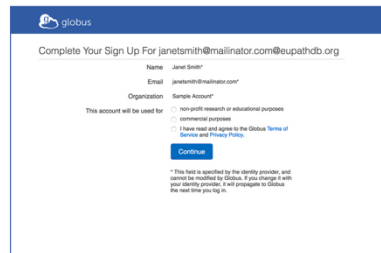
## Analyze My Experiment

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).

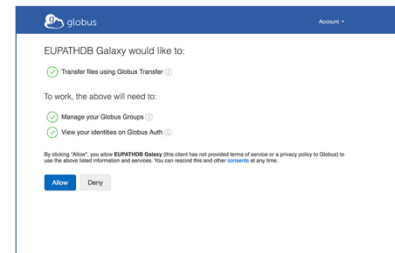Click **"Continue to Galaxy"** to sign up for EuPathDB Galaxy services.

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**.
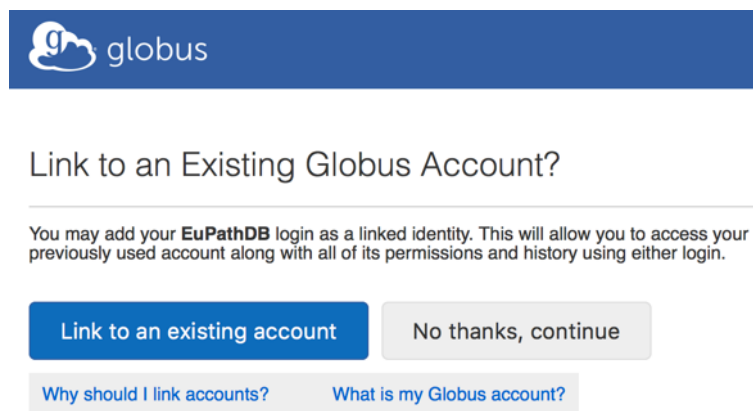
(3) Grant permission to share your Globus identity and files with us. Please click **Allow**. (We will only perform file transfers that you explicitly request, between Galaxy and other resources, including EuPathDB.)
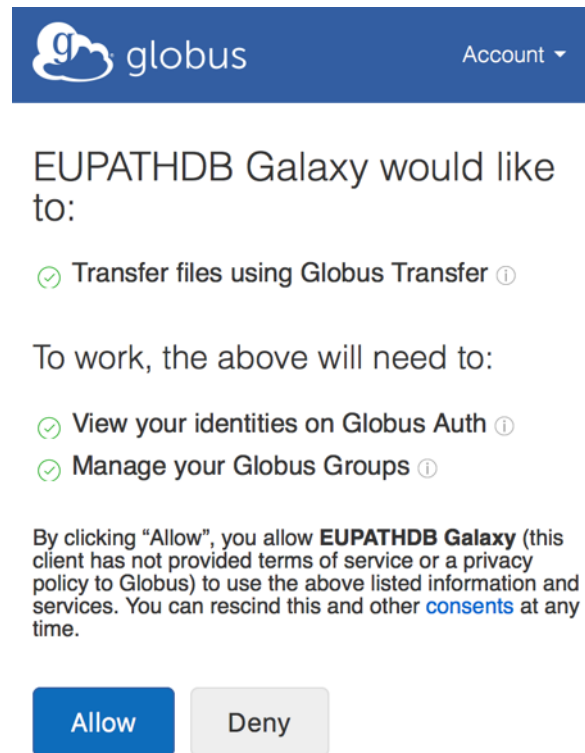
Continue to Galaxy

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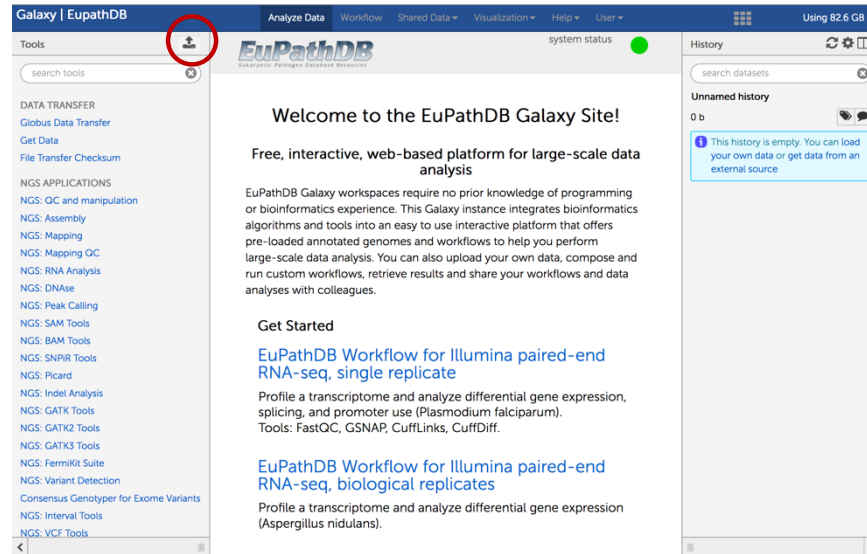Step 6: Click on "No thanks, continue"

Step 7: Click on "Allow"
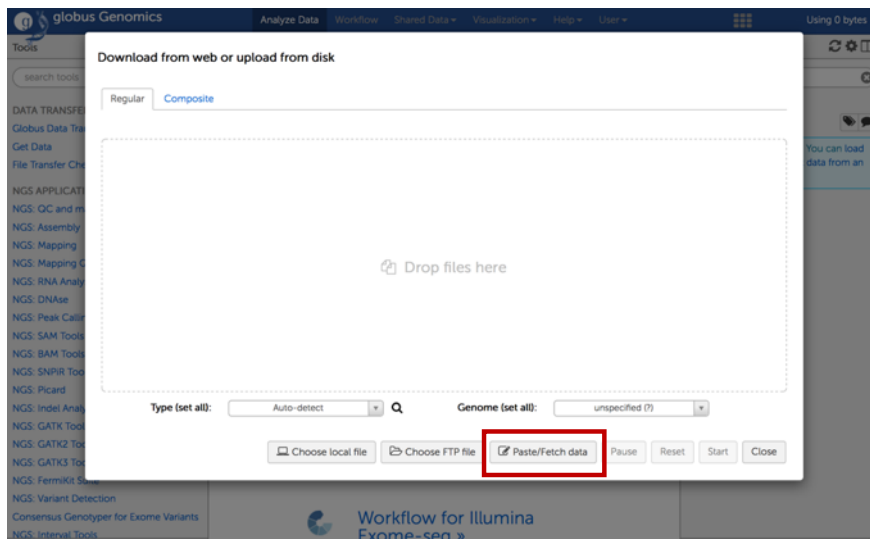


Step 8: Congratulations, you are in!

## Section II: Importing data to Galaxy

There are multiple ways to important data into your Galaxy workspace. For this exercise we will use the direct links listed below. Remember one person in your group will be doing this. 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" icon. This will open up a window that allows you to "**Download from web or upload from disk**"

Step 2: In the "Download from web or upload from disk" window click on "Paste/Fetch data"

Step 3: Paste the four URLs corresponding to the four files for your group. Each URL has to be on a new line. Then click on "Start".



Step 4: Click on "Close". You should notice that the left section (history section) will show the files being transferred (yellow) – this may take a few minutes to start. File transfer will take about 15-20 minutes. When this is complete they will turn green.

## Group assignments:

**Group 1:**

*Plasmodium falciparum* Asexual vs. Cultured sporozoites
Project information: http://www.ebi.ac.uk/ena/data/view/PRJNA230379

Samples:

Asexual samples:
ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR104/000/SRR1041270/SRR1041270_1.fastq.gz
ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR104/000/SRR1041270/SRR1041270_2.fastq.gz

Cultured sporozoite samples:
ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR104/008/SRR1041268/SRR1041268_1.fastq.gz
ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR104/008/SRR1041268/SRR1041268_2.fastq.gz

**Group 2:**

*Plasmodium falciparum* Asexual vs. Salivary sporozoites
Project information: http://www.ebi.ac.uk/ena/data/view/PRJNA230379

Samples:

Asexual samples:
ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR104/000/SRR1041270/SRR1041270_1.fastq.gz
ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR104/000/SRR1041270/SRR1041270_2.fastq.gz

Salivary Sporozoites:
ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR104/006/SRR1041266/SRR1041266_1.fastq.gz
ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR104/006/SRR1041266/SRR1041266_2.fastq.gz

**Group 3:**

*Plasmodium falciparum* Cultured vs. Salivary sporozoites
Project information: http://www.ebi.ac.uk/ena/data/view/PRJNA230379

Samples:

Cultured sporozoite samples:
ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR104/008/SRR1041268/SRR1041268_1.fastq.gz
ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR104/008/SRR1041268/SRR1041268_2.fastq.gz

Salivary Sporozoites:
ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR104/006/SRR1041266/SRR1041266_1.fastq.gz
ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR104/006/SRR1041266/SRR1041266_2.fastq.gz

**Group 4:**

*Aspergillus nidulans* FGSC4 VeA[+] WT vs. OSA knock outs
Project information: http://www.ebi.ac.uk/ena/data/view/PRJNA293709

Samples:

FGSC4 VeA[+] WT:
ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR218/001/SRR2180251/SRR2180251_1.fastq.gz
ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR218/001/SRR2180251/SRR2180251_2.fastq.gz

OSA knock outs:
ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR218/007/SRR2180257/SRR2180257_1.fastq.gz
ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR218/007/SRR2180257/SRR2180257_2.fastq.gz

**Group 5:**

*Toxoplasma gondii* WT vs. GRA17 knock outs
Project information: http://www.ebi.ac.uk/ena/data/view/PRJNA275621

Samples:

WT:
ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR180/001/SRR1805881/SRR1805881_1.fastq.gz
ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR180/001/SRR1805881/SRR1805881_2.fastq.gz

GRA17 knock outs:
ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR180/002/SRR1805882/SRR1805882_1.fastq.gz
ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR180/002/SRR1805882/SRR1805882_2.fastq.gz

**Group 6:**

*Toxoplasma gondii* WT vs. GRA17 knock outs
Project information: http://www.ebi.ac.uk/ena/data/view/PRJNA275621

Samples:

WT:
ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR180/001/SRR1805881/SRR1805881_1.fastq.gz
ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR180/001/SRR1805881/SRR1805881_2.fastq.gz
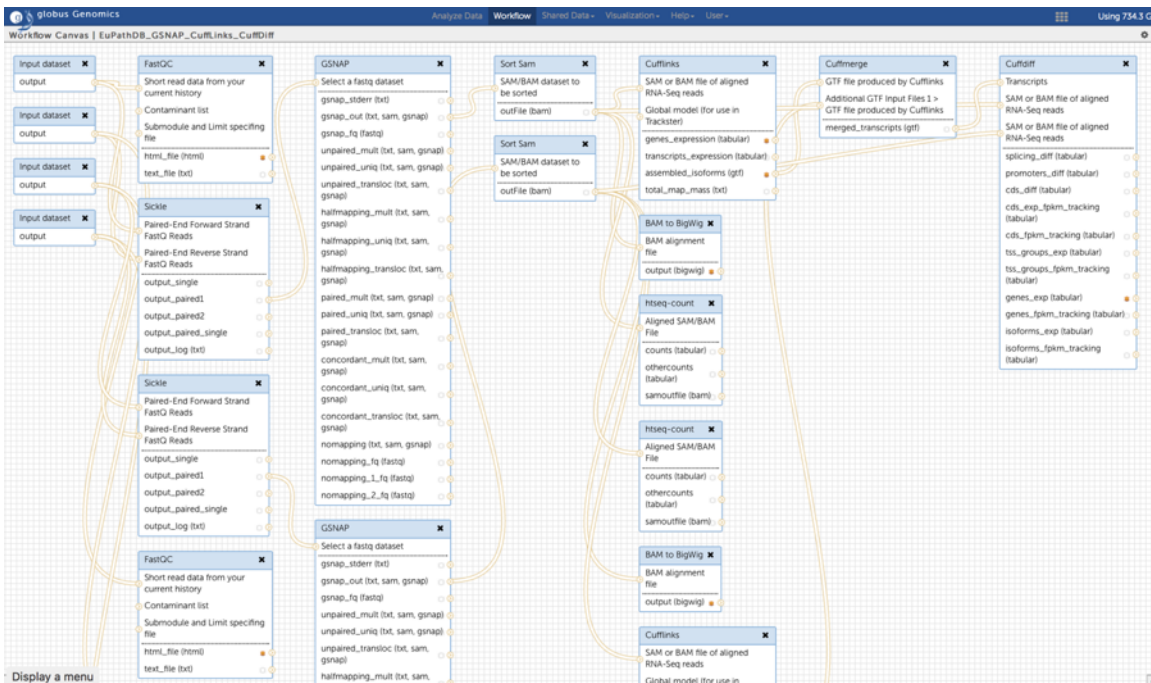
GRA23 knock outs:
ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR180/003/SRR1805883/SRR1805883_1.fastq.gz
ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR180/003/SRR1805883/SRR1805883_2.fastq.gz

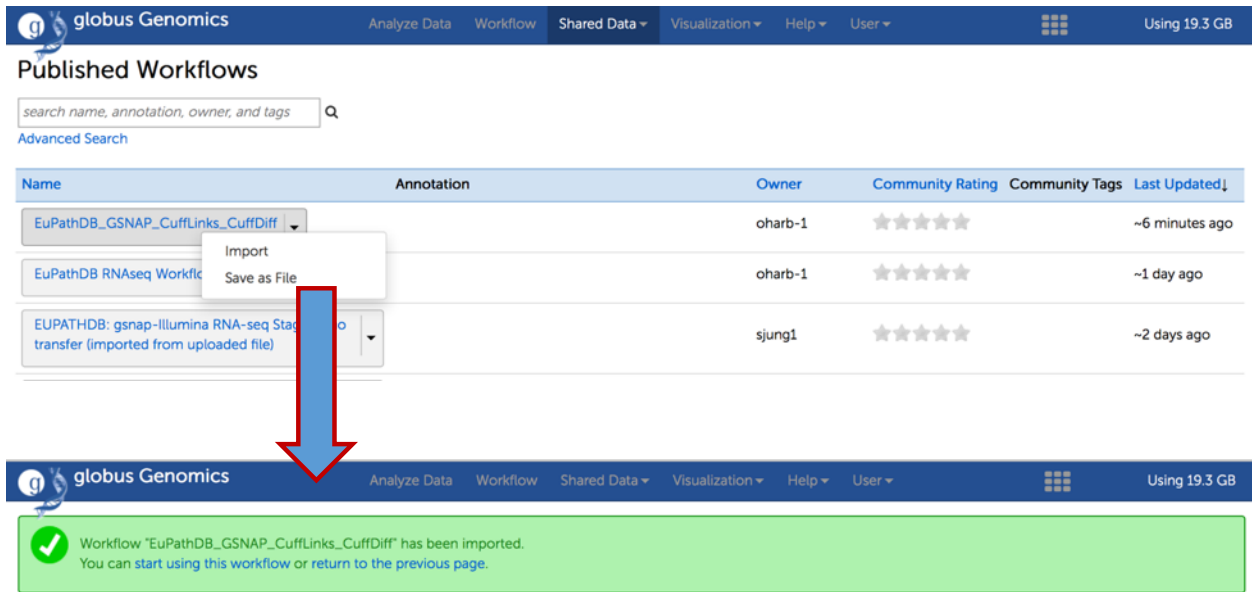## Section II: Running a workflow in Galaxy

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.
3. Aligns the reads to a reference genome using GSNAP and generates coverage plots.
4. Determines FPKM values for each sample and generates gene/transcript models.
5. Determines differential expression of genes between the samples.



Step 1: Import the workflow called "EuPathDB_GSNAP_CuffLinks_CuffDiff" – click on the shared data menu item and select "Published Workflows" from the menu.

**Step 2:** Click on the arrow next to the appropriate workflow and select import.



**Step 3:** Click on "Workflow" in the menu at the top of the page. On the next page click on the arrow next to your imported workflow and select the "Run" option.

Step 4:  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 datasets.  These are the fastq files you imported from the sequence archive.  Workflow steps 1-4 allow you to select the datasets.  Be sure you match the correct forward and reverse files.  The should end in the same SRR number with a .1 or .2 at the end.

b. Scroll down to steps 11 and 12 (GSNAP).  Click on the name of the step to open up the parameters.  Select the correct reference organism in each of the steps.

Step 11: GSNAP (version GSNAP: 2014-08-04)

7

&lt;H2&gt;Input Sequences&lt;/H2&gt;Select the input format

Fastq

Select a fastq dataset

Output dataset 'output_paired1' from step 6

Use Paired Reads?

False

Amount of barcode to remove from start of read (default 0)

None ✎

Starting field of identifier in FASTQ header, whitespace-delimited, starting from 1

None ✎

Ending field of identifier in FASTQ header, whitespace-delimited, starting from 1

None ✎

Skip reads marked by the Illumina chastity program

off - no filtering ✎

Select a reference genome

| AnidulansFGSCA4 | ⇕ ↻ |
|---|---|

TREU927 (Tbrucei)
hg19 (Hsapiens)
ME49 (Tgondii)
3D7 (Pfalciparum)
C57BL6J (Mmusculus)
PvivaxSal1
AfumigatusAf293
AnidulansFGSCA4

s

A-Seq

put options for RNA-Seq

Use default settings

c. Scroll down to step 15 (Cufflinks), 17, 18 (htseq), 20 (Cufflinks) and 21 (Cuffmerge) and select the correct reference organism.
d. Click on "Run Workflow"

Step 21: Cuffmerge (version CUFFLINKS: 2.1.1)

33

**GTF file produced by Cufflinks**
Output dataset 'assembled_isoforms' from step 15

**Additional GTF Input Files**

> **Additional GTF Input Files 1**
>
> **GTF file produced by Cufflinks**
> Output dataset 'assembled_isoforms' from step 20

**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**

| Pfalciparum 3D7 | ⇕ ↻ |
|---|---|

**Use Sequence Data**
No

**Action:**
Hide output 'merged_transcripts'.

Step 22: Cuffdiff (version CUFFLINKS: 2.1.1)

23

☐ Send results to a new history

Run workflow

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