RNA sequence data analysis via VectorBase Galaxy

Part I: Uploading data and starting the workflow

Learning objectives

- Become familiar with VectorBase Galaxy workspace
- Import data from EBI to the VectorBase Galaxy
- Create collections of datasets
- Run a pre-configured RNA-Seq workflow

The VectorBase Galaxy-based workspace offers pre-loaded genomes, private data analysis and display, and the ability to share and export analysis results and also import certain datasets into private workspace within VectorBase (*My Workspace* > My data sets).

VectorBase Galaxy workspace can be accessed from the *My Workspace* tab or from the *Tools* menu tab, on the home page of VectorBase or any other VEuPathDB site. To log in, users must have an account with VectorBase/VEuPathDB, which is open access. After an account is created, users receive access to the VectorBase Galaxy services and tools.

| VEuPathDB ^{Release 52} 20 May 2021 | Site search, e.g. PF3D7_1133400 | Site search, e.g. PF3D7_1133400 or *reductase or "binding protein" | | | | |
|--|---------------------------------|--|--|--|--|--|
| Eukaryotic Pathogen, Vector & Host Informatics Res | My Strategies Searches Tools | My Workspace Data About Help Contact Us Analyze my data (Galaxy) nderstand how you | | | | |
| Improve the service. Your opinion matters! https: Search for | //bit.ly/33Gpo51 | My baskets My BLAST beta jobs | | | | |
| expand all I collapse all | | Public search strategies | | | | |

The Galaxy instance is not meant for long-term data storage. Datasets are automatically deleted after 60 days or when the total quota for all projects is reached. To save your data, download your analysis results locally and then *delete and purge* files to free up space for your next analysis; to do it, follow these steps.

Gear icon > History options



History lists > Saved Histories

| | Analyze Data | | | | | | | Using 188.4 GB |
|--------------------|--|--|---|---|---------------------------------------|---------------------------------|-------------------------|--|
| Atte not kee | VEUPathDB Letryale Pelages. Tester 4 Best ention: We period been updated in t p long term, and c | internatio Reserves lically purge y he past 60 da delete unwant | our inactive datası ıys will be removec ed histories. Thanl | ets (inputs and out 1. Please backup/dı k you for your coop | puts). Any ownload yo peration. | GG-v dataset t our datase | 5.4 hat has ts to | History LISTS Saved Histories An Histories Shared with Me PRI CURRENT HISTORY |

Click on the history of interest (down arrowhead) and "Delete Permanently".

| Sa | Saved Histories | | | | | |
|--------|-----------------------------|----------------------------------|--|--|--|--|
| sear | ch history name | es and tags Q | | | | |
| Adva | nced Search | | | | | |
| | Name | Datasets | | | | |
| | (unstranded) 28:unx vs 4 | PERU 2 Switch | | | | |
| | (unstranded 35:acyp vs | View Share or Publish Copy | | | | |
| | (unstranded 14:unx vs 2 | Rename Delete | | | | |
| \Box | (unstrand | Delete Permanently | | | | |

Optionally, you can also transfer your data for example to a cloud service with Data transfer > Globus Data Transfer



Part I: Uploading data and starting the workflow

For this exercise (Part I), 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. Set the workflow to run overnight and view / interpret the results the next day (Part II).

Setting up your VectorBase Galaxy account

The step-by-step process has been explained in the other companion handout of this webinar.

The anatomy of the VectorBase Galaxy landing page

The workspace has four major components:

a) the top menu controls the main interface

- b) the left panel has a list of available tools
- c) the main welcome page is the interactive interface that houses pre-configured workflows, workflows editor, etc.
- d) the right panel provides access to histories, deleted datasets, and other useful functions



Section II: Importing data to Galaxy

There are multiple ways to important data into your Galaxy workspace. For this exercise, we will use the 'Get Data via Globus from the EBI: server using your unique file identifier" tool and enter the sequence repository sample IDs. The samples below were all generated by paired end sequencing; hence each sample ID will result in transferring two files to your galaxy history. The files are fastq files that are compressed (that is why they end in .gz = gzip).

Sample data set:

We will be examining data from a study called "Transcriptome analysis of genes associated with pyrethroid resistance in South and Central American *Anopheles albimanus*". <u>https://www.ncbi.nlm.nih.gov/bioproject/PRJNA498101</u>.

| Sample: Accession Numbers (Name) | Sanarate | Guatemala | Peru |
|--|---|--|--|
| unexposed to insecticides | SAMN10341948 (San3) SAMN10341947 (San2) SAMN10341946 (San1) | SAMN10341945 (GTM_unx3) SAMN10341944 (GTM_unx2) SAMN10341943 (GTM_unx1) | SAMN10341939 (PER_unx3) SAMN10341938 (PER_unx2) SAMN10341937 (PER_unx1) |
| delta-methrin | - | SAMN10341942 (GTM_delta3) SAMN10341941 (GTM_delta2) SAMN10341940 (GTM_delta1) | SAMN10341933 (PER_delta3) SAMN10341932 (PER_delta2) SAMN10341931 (PER_delta1) |
| alpha- cypermethrin | - | - | SAMN10341936 (PER_acyp3) SAMN10341935 (PER_acyp2) SAMN10341934 (PER_acyp1) |

Reference Genome in VectorBase: Anopheles albimanus STECLA AalbS2.7.

Comparisons: We will compare the San to the GTM-delta samples, you can later do the other comparisons following the same steps.

| | Condition- 1 | Condition-2 | |
|---|-----------------|-------------|------------------|
| < | San | GTM-delta | \triangleright |
| | San | GTM-unx | |
| | GTM-unx | GTM-delta | |
| | San | PER-delta | |
| | San | PER-unx | |
| | San | PER-acyp | |
| | PER-unx | PER-delta | |
| | PER-unx | PER-acyp | |
| | PER-delta | PER-acyp | |
| | GTM-delta | PER-delta | |

Step 1: Click on the "**Globus Data Transfer**" link in the left-hand menu. This will reveal a list of options; click on "**Get Data via Globus from the EBI server**". ***important: do not select the option for transferring a collection.

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 *Anopheles albimanus* susceptible laboratory colony Sanarate, San (three biological replicates, San1, San2, San3) and field collected mosquitoes from Guatemala and Peru (three biological replicates each, unexposed and exposed to two different insecticides), it is useful to organize these into collections. Using collections also reduces the complexity of the Galaxy workflows. See below:





| 42: Peru alive after exposure to deltamethrin (PER-delta) a list of 3 dataset pairs | × |
|---|---|
| 35: Peru alive after exposure to alpha-cypermethrin (PER-acyp) a list of 3 dataset pairs | × |
| 28: Peru unexposed to insecticid e (PER-unx) a list of 3 dataset pairs | × |
| 21: Guatemala alive after expos ure to deltamethrin (GTM-delta) a list of 3 dataset pairs | × |
| 14: Guatemala unexposed to ins ecticide (GTM-unx) a list of 3 dataset pairs | × |
| 7: Sanarate susceptible lab colo ny (San) a list of 3 dataset pairs | × |

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).



Additional resources:

Galaxy Project: <u>https://usegalaxy.org/</u> Trimmomatic: <u>http://www.usadellab.org/cms/index.php?page=trimmomatic</u> FastQC: <u>https://www.bioinformatics.babraham.ac.uk/projects/fastqc/</u> HISAT2: <u>http://daehwankimlab.github.io/hisat2/main/</u> HTseq: <u>https://htseq.readthedocs.io/en/master/</u> DEseq2: <u>https://doi.org/10.1186/s13059-014-0550-8</u>

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

| 🕡 👌 globus Genomics | Analyze Data Workflow Shared Data - Visualization - Help - User - |
|--------------------------|---|
| Tools | Get started with VEuPathDB pre-configured workflows: |
| search tools | |
| | OrthoMCL |
| VEUPATHDB APPLICATIONS | This workflow uses BLASTP and the OrthoMCL algorithm to assign your set of proteins to OrthoMCL groups. Usually, the |
| VEUPathDB Export Tools | most recent version (e.g., OG6r5) is in sync with the groups on the OrthoMCL website, but you can also select a previous set such as OG5. Explore this OrthoMCL workflow tutorial to learn more. |
| /EuPathDB OrthoMCL Tools | Mortflew to man your proteins to OrthoMCL groups |
| Vedratibb Kitk Seq 1003 | • Workhow to map your proteins to orthonice groups |
| DATA TRANSFER | RNA-seq |
| Globus Data Transfer | Use the following workflows to analyze your FASTQ files. The workflows use FASTQ groomer and Trimmomatic for |
| Get Data | preparation of reads, FASTQC for sequencing statistics, and HISAT2 for mapping reads to a VEuPathDB reference genome |
| Collection Tools | Choose the appropriate workflow based on your input data and your desired analysis. Explore this RNA-Seq export tutorial to learn about exporting your workflow results to VEN24tbDB. |
| ile Transfer Checksum | catorial to learn about exporting your worknow results to vicuration. |
| NGS APPLICATIONS | Examine genome coverage and calculate TPM for each gene |
| NGS: QC and manipulation | In addition to the tools described above, these workflows use three tools (bamCoverage, htseq-count, HTSeqCountTOTPM to generate BinWig and TPM files that can be analyzed on VENPathBR in Galaxy or on your computer. The workflows tak |
| IGS: Assembly | any number of samples and processes them in parallel. To export the results to VEUPathDB, use the 'RNA-Seq to |
| IGS: Mapping | VEuPathDB' tool. |
| GS: Mapping QC | Workflow for paired-end stranded reads |
| IGS: HLA Typing | Workflow for paired-end unstranded reads Workflow for inclusion detranded reads |
| IGS: RNA Analysis | Workflow for single-end sublided reads Workflow for single-end unstranded reads |
| IGS: miRNA | Identify genes with statistically significant expression differences between two samples |
| IGS: DNAse | In addition to the tools described above, these workflows use three tools (htseq-count, DESeq2, Bam to BigWig) to |
| YS: MTB | determine whether each gene exhibits differential expression and to generate BigWig coverage files. The output files can |
| GS: Mothur | be analyzed in Galaxy or on your computer. The worknows compare two samples with any number of replicates. To expor your BioWin files to VFUPathDR, use the 'Biowin Files to VFUPathDR' tool. To filter your DESeo2 result file and obtain a set |
| GS: QIIME | , Gene IDs that change significantly (defaults: fold-change>=2 and adj-p<=0.05; these can be changed), use this |
| GS: QIIME2 | workflow. Copy and paste the Gene IDs into the 'Identify Genes based on Gene ID(s)' question on a VEuPathDB website, |
| GS: PICRUST | as seen nere for the Plasmous site. |
| GS: Parallel-Meta | Workflow for paired-end stranded reads Workflow for paired-end unstranded reads |
| GS: BIOM | Workflow for single-end stranded reads |
| IGS: DADA2 | Workflow for single-end unstranded reads |
| NGS: HOMER | Variant calling |
| IGS: Peak Calling | Variant Canny |

- 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.

| Workflow: imported: RNA-Seq DESeq2 PE unstranded (v.4) | ✓ Run workflow |
|--|----------------|
| <u>1: Input Dataset Collection - Sample 1</u> | - |
| 7: Sanarate susceptible lab colony (San) | |
| 2: Input Dataset Collection - Sample 2 | |
| 21: Guatemala alive after exposure to deltamethrin (GTM-delta) | • • |

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 that each of these tools is in the workflow twice since you have two collections). It is critical that you select the correct genome that matches the experimental organism. For this exercise, select the genome as shown in the image below

| Select a reference genome | VectorBase-49_AaegyptiLVP_AGWG_Genome |
|--|---|
| VectorBase-49_AalbimanusSTECLA_Genome | VectorBase-49_AalbimanusSTECLA_Genome |
| If your genome of interest is not listed, contact the Galaxy team Primary alignments | VectorBase-49_AarabiensisDongola_Genome |

c. Two additional parameters to check in the htseq-count step are "Feature type" and "ID Attribute". They should be set to "exon" and "gene_id", respectively. Be aware that these are case-sensitive, so "Exon" is not correct but "exon" is correct. Here is how that step should look:

| · · · |
|---|
| 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 10 |
| 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 |
| VectorBase-49_AalbimanusSTECLA_Genome |
| |
| Union |
| ☑ Stranded |
| No |
| 🕼 Minimum alignment quality |
| 10 |
| 🕼 Feature type |
| exon |
| 🕼 ID Attribute |
| gene_id |

n.

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

| Workflow: imported: RNA-Seq DESeq2 PE unstranded (v.4) | | | | | |
|--|--|--|--|--|--|
| History Options | | | | | |
| Send results to a new history | | | | | |
| Yes No | | | | | |
| <u>1: Input Dataset Collection - Sample 1</u> | | | | | |
| ☐ ■ 7: Sanarate susceptible lab colony (San) | | | | | |
| 2: Input Dataset Collection - Sample 2 | | | | | |
| 21: Guatemala alive after exposure to deltamethrin (GTM-delta) | | | | | |

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.

| | History | C 🌣 🗆 |
|--|--|----------|
| Successfully invoked workflow imported: RNA-Seq DESeq2 PE unstranded (v.4) . You can check the status of queued jobs and view the resulting data by refreshing the | search datasets | 8 |
| completed successfully or 'error' if problems were encountered. | (unstranded) San vs G 6 shown, 36 hidden | TM-delta |
| | 18.59 GB | 2 🃎 🗩 |

Practice working with Galaxy editor

You can create your own workflows. The tools can all be added and configured in an interactive workflow editor.

• Navigate to the Workflow tab from the main menu at the top



• Left click on the drop-down icon within the workflow you want to modify and select the "Edit" option.



• Delete HISAT2 step by clicking on the "x" in the top right corner.

• Locate the HISAT2 tool in the Tools panel and click to insert it back into the workflow.



- Re-establish connections for HISAT2
- Click on the arrow in the step before HISAT2 and drag to the appropriate input in HISAT2 tool.
- What happens? Can you reconnect it?



Note: Sometimes you may be unable to re-establish connection. When this happens, take a look at the tool documentation notes in the right panel, check if your selection for single-read or paired-end setting in particular (paired-end setting must be selected if you are dealing with reverse and forward reads).

| Details | |
|-----------------|---------------------------|
| Save | AT2 A fast and |
| Save As | ive alignment |
| Run | am (Galaxy Version 2.0.5) |
| Edit Attributes | |
| Auto Re-layout | |
| Close | |
| tep label. | |
| Annotation | |

Now that you have learned the principals of workflow editing, you can either practice saving the workflow by clicking on the wheel at the far top corner or simply existing the workflow editor without saving.

Appendix:

FASTQ files 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 (SRA) at GenBank
- The European Nucleotide Archive (ENA) at EMBL
- The DNA Data Bank of Japan

