

INF-BIO5121/9121 - Oct 7, 2014

Analyzing miRNA data using Lifeportal

PRACTICALS

In this experiment we have miRNA data from the livers of baboons (*Papio Hamadryas*) before and after they are given a high cholesterol high fat diet. We have three samples from baboons on the baseline diet (Chow), and three samples from the baboons on the high cholesterol high fat diet (HCHF). The goal of the analysis is to find miRNAs that are differentially expressed in these two groups. We will start out with the fastq read files and follow these steps to get the results:

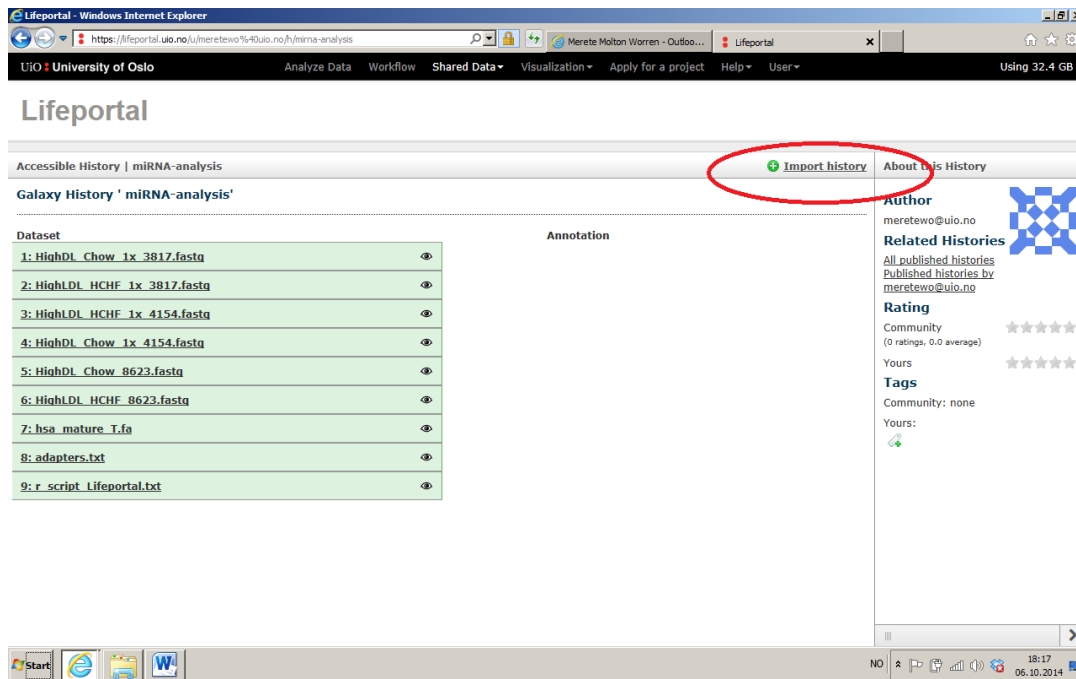
1. Logging in to lifeportal and get the read files from a shared history
2. Removing adapters
3. Looking at read qualities with FastQC
4. Align the reads to human miRNAs
5. Count how many reads map to each miRNA
6. Make a workflow, and run it
7. Use join and cut to get the appropriate file format
8. Run a differential expression analysis with EdgeR
9. Make a Galaxy page

1. Getting the read files

Start by using this link to go to the history page:

<https://lifeportal.uio.no/u/sveinugu%40uio.no/h/infbio121-mirna-demo>

You should then see a page like this:



S

Click **import history**, and on the next page click **start using this history**.

This will create a new history for you to use, with all the needed files in it.

In the right pane, you will have your history items. In each of them you will find a little symbol of an eye. Click on this to view the data, for at least one of the files. Expand one of the history items by clicking on the name, and see what kind of info you get there.

2. Removing adapters

Since miRNAs are shorter (~22 bp) than the read length (36 bp), the reads will contain a stretch of the 3' adapter. These should be removed before further analysis is done. In Lifeportal, under **NGS:QC and manipulation** we can find the **fastx-toolkit** with an option to **Clip adapter sequences** that can be used for this. Below is a picture of how it looks after clicking the **Clip**. Take a look at it before starting.

There are several options that needs to be set before running this program.

First, choose a file for the **Library to clip** option. A list will come up where you can choose from the ones Lifeportal can see has the right format. Choose the first of the provided HCHF fastq read files (("1: HighLDL_HCHF_1x_3817").

Next, set the **Minimum sequence length** to 17.

Then it is time to put in the adapter sequence that is to be removed. Under **Source** we choose **Enter custom sequence**. The 3' adapter sequence is ATCTCGTATGCCGTCTTCTGCTTGT, copy and paste this sequence into the **Enter custom clipping sequence** field.


In the next field, leave the 0. We do not want to keep the adapter sequence.

For **Discard sequences with unknown (N) bases**, choose **Yes**

For **Output options** choose **Output only clipped sequences**

Now, all you need to do is to choose a project under **Projects/Accounts** and click the blue **Execute** button.

While your job is running, click the **back arrow** on your browser and check the description of what it does at the bottom of the middle pane.

When your job is finished, look at the file using the little eye button in your history item. You should see right away that the reads are shorter than they were before the clipping. Click on the history item. After the description field you will find a **rerun** button . Click on it to get the options for this run back in the middle pane. That way, for the next **Clip adapter** run there will be less work.

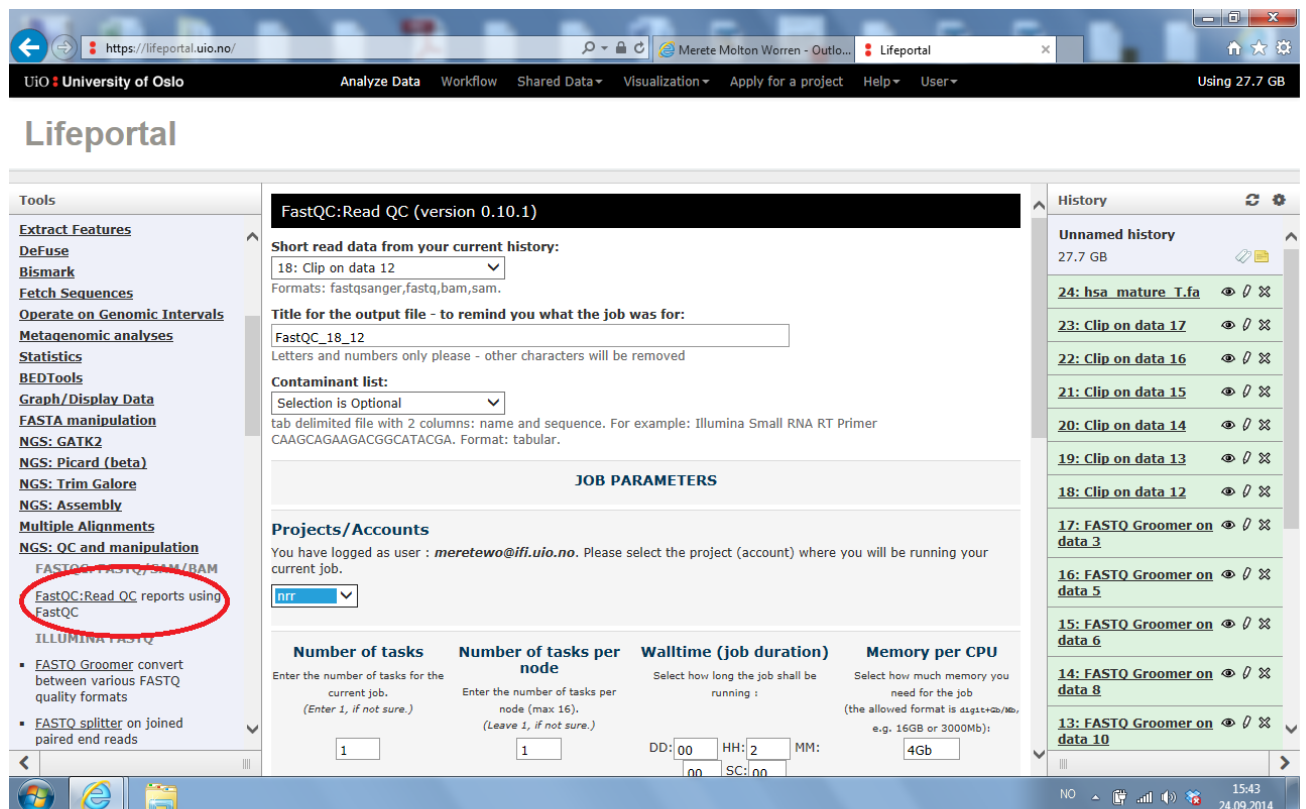
For now, we will only work with this file, so do not rerun the job.

3. Looking at read qualities with FastQC

Now that the adapter is removed, it is time to look at the reads with FastQC, to make sure they look ok.

Choose **FastQC:Read QC report using FastQC** under **NGS:QC and manipulation**.

Choose the Clip history item you just made. Make a title for the output. We do not use a contamination list in this instance. Choose the project and click execute.



The screenshot shows the Lifeportal web interface. The browser address bar displays <https://lifeportal.uio.no/>. The top navigation bar includes links for **Analyze Data**, **Workflow**, **Shared Data**, **Visualization**, **Apply for a project**, **Help**, and **User**. The main content area is titled **Lifeportal** and features a sidebar on the left with a **Tools** menu. The **Tools** menu is expanded, showing categories like **Extract Features**, **DeFuse**, **Bismark**, **Fetch Sequences**, **Operate on Genomic Intervals**, **Metagenomic analyses**, **Statistics**, **BEDTools**, **Graph/Display Data**, **FASTA manipulation**, **NGS: GATK2**, **NGS: Picard**, **NGS: Trim Galore**, **NGS: Assembly**, **Multiple Alignments**, and **NGS: QC and manipulation**. The **NGS: QC and manipulation** category is selected, and the **FastQC:Read QC reports using FastQC** tool is highlighted with a red circle. The main panel displays the configuration for **FastQC:Read QC (version 0.10.1)**. It includes a dropdown for **Short read data from your current history:** set to **18: Clip on data 12**. The **Title for the output file** is **FastQC_18_12**. The **Contaminant list** is set to **Selection is Optional**. The **JOBS PARAMETERS** section shows **Projects/Accounts** with a dropdown set to **lrr**. The **Number of tasks** is **1**, **Number of tasks per node** is **1**, **Walltime (job duration)** is **DD:00 HH:2 MM:00**, and **Memory per CPU** is **4Gb**. The right sidebar shows the **History** panel with a list of jobs, including **24: hsa_mature_T.fa**, **23: Clip on data 17**, **22: Clip on data 16**, **21: Clip on data 15**, **20: Clip on data 14**, **19: Clip on data 13**, **18: Clip on data 12**, **17: FASTQ Groomer on data 3**, **16: FASTQ Groomer on data 5**, **15: FASTQ Groomer on data 6**, **14: FASTQ Groomer on data 8**, and **13: FASTQ Groomer on data 10**. The bottom status bar shows the system time as **15:43** on **24.09.2014**.

4. Mapping

Now that the reads are cleaned up, it is time to map them to the reference. The reference we use are the human mature miRNAs found in miRBase. A file with the sequences are found in your imported history ("7: hsa_mature_T.fa").

Start by finding **Bowtie2** under **NGS:Mapping** and click on it.

Choose **Single-end** under **Is this library mate paired?**.

Choose the result of the previous Clip job as the Fastq file.

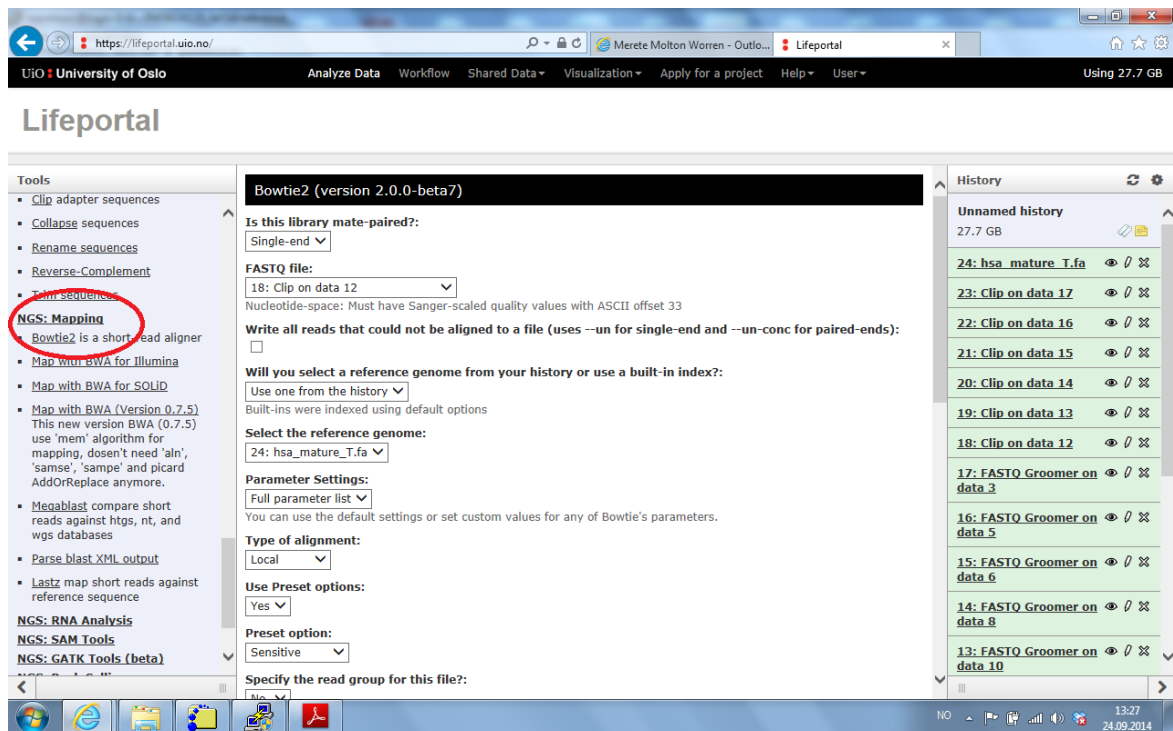
Do not check the box for writing unaligned reads to a file.

Choose **Use one from the history** under the question **Will you select a reference genome from your history or use a built in index?**.

Select the history item called **hsa_mature_T.fa** as your reference.

Choose to have the **full parameter list**. Change the **Type of alignment** to **Local**, and leave the others as they are.

Choose your project and click execute.



When the run is finished, the result is a .bam file, which is a binary file. Since we want to count the occurrences of the different miRNAs using textual manipulation, we need to transform it to a .sam file, which is a tabular file (i.e. a tab-delimited text file).

Under **NGS:SAM tools** choose **BAM-to-SAM**. Choose the Bowtie2 result file.

Do not check **Include header in output**.

Choose your project and run.

The screenshot shows the Lifeportal web interface at <https://lifeportal.uio.no/>. The user is logged in as Merete Molton Worren. The main navigation bar includes links for Analyze Data, Workflow, Shared Data, Visualization, Apply for a project, Help, and User. The page title is "Lifeportal".

The left sidebar lists various tools under the "Tools" section. The "NGS: SAM Tools" category is expanded, and the "BAM-to-SAM" tool is highlighted with a red circle. Other tools listed include Filter SAM, Convert SAM to interval, SAM-to-BAM, Merge BAM files, MPileup, Generate pileup, Filter pileup, Pileup-to-Interval, flagstat, rmdup, and Slice BAM.

The main content area displays the "BAM-to-SAM (version 0.1.18)" tool configuration. The "BAM File to Convert" dropdown is set to "32: Bowtie2 on data 18 and data 24: aligned reads". The "Include header in output" checkbox is unchecked. Below this is the "JOB PARAMETERS" section, which includes a "Projects/Accounts" dropdown set to "int".

The "JOB PARAMETERS" section contains four input fields:

- Number of tasks:** Enter the number of tasks for the current job. (Enter 1, if not sure.) Value: 1
- Number of tasks per node:** Enter the number of tasks per node (max 16). (Leave 1, if not sure.) Value: 1
- Walltime (job duration):** Select how long the job shall be running. DD: 00, HH: 2, MM: 00, SC: 00
- Memory per CPU:** Select how much memory you need for the job (the allowed format is digit+G/Mb, e.g. 16GB or 3000Mb). Value: 4Gb

An "Execute" button is located below the input fields. The bottom section, "What it does", is partially visible.

The right sidebar shows a "History" section with a list of previous jobs, including "33: Bowtie2 on data 19 and data 24: aligned reads", "32: Bowtie2 on data 18 and data 24: aligned reads", "31: FastQC2317 Clip on data 17.html", "30: FastQC2216 Clip on data 16.html", "29: FastQC2115 Clip on data 15.html", "28: FastQC Clip on data 14.html", "27: adapters.txt", "26: FastQC1913 Clip on data 13.html", "25: FastQC1812 Clip on data 12.html", "24: hsa_mature_T.fa", "23: Clip on data 17", and "22: Clip on data 16".

5. Counting

The next thing we need to do is to count the reads for each individual miRNA.

Under **Statistics**, choose **Count occurrences of each record**.

Under **from dataset** choose your **BAM-to-SAM** history item.

Count occurrences of values in column(s) should be c3.

Keep the **delimited by Tab** and run.

The screenshot shows the Lifeportal web interface. The top navigation bar includes 'Uio University of Oslo', 'Analyze Data', 'Workflow', 'Shared Data', 'Visualization', 'Apply for a project', 'Help', and 'User'. The main header is 'Lifeportal'. On the left, a 'Tools' sidebar lists various analysis tools, with 'Count occurrences of each record' highlighted under the 'Statistics' section. The main panel displays the 'Count (version 1.0.0)' tool configuration. It shows the 'from dataset' dropdown set to '12: BAM-to-SAM on data 11: converted SAM'. Below this, 'Count occurrences of values in column(s):' is set to 'c3'. The 'Delimited by' dropdown is set to 'Tab'. The 'JOB PARAMETERS' section includes 'Projects/Accounts' (set to 'nrr') and four input fields: 'Number of tasks' (1), 'Number of tasks per node' (1), 'Walltime (job duration)' (DD: 00, HH: 2, MM: 00, SC: 00), and 'Memory per CPU' (4Gb). An 'Execute' button is at the bottom. A tip at the bottom states: 'TIP: If your data is not TAB delimited, use Text Manipulation->Convert'. On the right, a 'History' panel shows a list of previous jobs, including 'imported: INFBI0x121-miRNA-demo' and several 'HighDL' jobs.


Click the **eye** icon of the resulting history item. Note that the first column in the file contains the counts and the second column contains the name of the miRNA that is counted. In order to better manipulate the data, we need to reverse the order of the two columns. To do this, we will use the **cut** tool.

Under **Text Manipulation**, select the **cut** tool. This is typically used for extracting some of the columns, in any order. In our case, we will use the tool to extract *both* columns, but in the reverse order, i.e. column 2 before column 1.

Write **c2,c1** in the **cut columns** box.


Select the result of the **Count** job after **From**. Select a project and execute the job.

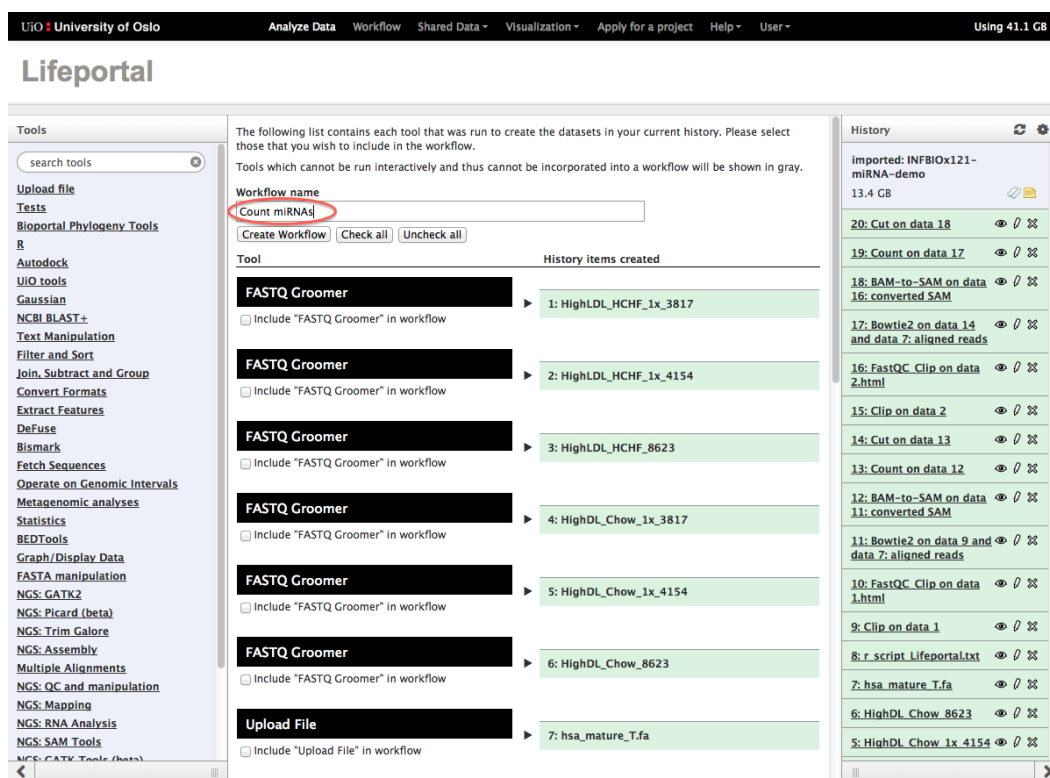
2-5. NOW, REDO ALL THE STEPS FOR THE NEXT HCHF FASTQ FILE (2: HighLDL_HCHF_1x_4154).

You should redo all six steps, from **Clip** to **Cut**. To make this easier, you can expand each of the previous history items and click the **rerun** button . Then change the input dataset accordingly, but keep the other options.

Note that you do not need to wait for each step to finish before you rerun the next one. Each step can be set up in advance, and will run when the previous step is finished. Note also that if you rerun the steps in this fashion (without waiting for each step to finish), you need to manually write “c3” in the **Count** step instead of selecting c3 from a list.

6. Create a workflow and rerun the steps for the last HCHF file

Click the **cogwheel** icon  on the top left of the history pane, and select **Extract workflow**. In the middle pane, you will now be able to select which previous jobs to base the new workflow on.



The screenshot displays the Lifeportal web interface. The top navigation bar includes links for 'Analyze Data', 'Workflow', 'Shared Data', 'Visualization', 'Apply for a project', 'Help', and 'User'. The main content area is divided into three panels:

- Tools:** A sidebar on the left lists various tools such as 'Upload file', 'Tests', 'Bioportal Phylogeny Tools', 'R', 'Autodock', 'UJO tools', 'Gaussian', 'NCBI BLAST+', 'Text Manipulation', 'Filter and Sort', 'Join, Subtract and Group', 'Convert Formats', 'Extract Features', 'DeFuse', 'Bismark', 'Fetch Sequences', 'Operate on Genomic Intervals', 'Metagenomic analyses', 'Statistics', 'BEDTools', 'Graph/Display Data', 'FASTA manipulation', 'NGS: GATK2', 'NGS: Picard (beta)', 'NGS: Trim Galore', 'NGS: Assembly', 'Multiple Alignments', 'NGS: QC and manipulation', 'NGS: Mapping', 'NGS: RNA Analysis', 'NGS: SAM Tools', and 'NGS: CATX Tools (beta)'.
- Workflow Creation:** The central panel shows a form to create a new workflow. The 'Workflow name' field is set to 'Count miRNAs'. Below this, there are buttons for 'Create Workflow', 'Check all', and 'Uncheck all'. A table lists tools and their associated history items:

Tool	History items created
FASTQ Groomer	1: HighLDL_HCHF_1x_3817
FASTQ Groomer	2: HighLDL_HCHF_1x_4154
FASTQ Groomer	3: HighLDL_HCHF_8623
FASTQ Groomer	4: HighDL_Chow_1x_3817
FASTQ Groomer	5: HighDL_Chow_1x_4154
FASTQ Groomer	6: HighDL_Chow_8623
Upload File	7: hsa_mature_T.fa
- History:** The right panel shows a list of previous workflow steps, including 'Imported: INFBIx121-miRNA-demo', '20: Cut on data 18', '19: Count on data 17', '18: BAM-to-SAM on data 16: converted SAM', '17: Bowtie2 on data 14 and data 7: aligned reads', '16: FastQC Clip on data 2.html', '15: Clip on data 2', '14: Cut on data 13', '13: Count on data 12', '12: BAM-to-SAM on data 11: converted SAM', '11: Bowtie2 on data 9 and data 7: aligned reads', '10: FastQC Clip on data 1.html', '9: Clip on data 1', '8: r_script Lifeportal.txt', '7: hsa_mature_T.fa', '6: HighDL_Chow_8623', and '5: HighDL_Chow_1x_4154'.

UiO University of Oslo Analyze Data Workflow Shared Data Visualization Apply for a project Help User Using 41.1 GB

Lifeportal

Tools
search tools
Upload file
Tests
Biportal Phylogeny Tools
R
Autodock
UiO tools
Gaussian
NCBI BLAST+
Text Manipulation
Filter and Sort
Join, Subtract and Group
Convert Formats
Extract Features
DeFuse
Bismark
Fetch Sequences
Operate on Genomic Intervals
Metagenomic analyses
Statistics
BEDTools
Graph/Display Data
FASTA manipulation
NGS: GATK2
NGS: Picard (beta)
NGS: Trim Galore
NGS: Assembly
Multiple Alignments
NGS: QC and manipulation
NGS: Mapping
NGS: RNA Analysis
NGS: SAM Tools
NGS: CATK Tools (beta)

☐ Include "Upload File" in workflow

Upload File
☐ Include "Upload File" in workflow
8: r_script_Lifeportal.txt

☒ Include "Clip" in workflow

Clip
☒ Include "Clip" in workflow
9: Clip on data 1

☒ Include "FastQC:Read QC" in workflow

FastQC:Read QC
☒ Include "FastQC:Read QC" in workflow
10: FastQC_Clip on data 1.html

☒ Include "Bowtie2" in workflow

Bowtie2
☒ Include "Bowtie2" in workflow
11: Bowtie2 on data 9 and data 7: aligned reads

☒ Include "BAM-to-SAM" in workflow

BAM-to-SAM
☒ Include "BAM-to-SAM" in workflow
12: BAM-to-SAM on data 11: converted SAM

☒ Include "Count" in workflow

Count
☒ Include "Count" in workflow
13: Count on data 12

☒ Include "Cut" in workflow

Cut
☒ Include "Cut" in workflow
14: Cut on data 13

☐ Include "Clip" in workflow

Clip
☐ Include "Clip" in workflow
15: Clip on data 2

☐ Include "FastQC:Read QC" in workflow

FastQC:Read QC
☐ Include "FastQC:Read QC" in workflow
16: FastQC_Clip on data 2.html

History
imported: INF8IOx121-miRNA-demo 13.4 GB
20: Cut on data 18
19: Count on data 17
18: BAM-to-SAM on data 16: converted SAM
17: Bowtie2 on data 14 and data 7: aligned reads
16: FastQC_Clip on data 2.html
15: Clip on data 2
14: Cut on data 13
13: Count on data 12
12: BAM-to-SAM on data 11: converted SAM
11: Bowtie2 on data 9 and data 7: aligned reads
10: FastQC_Clip on data 1.html
9: Clip on data 1
8: r_script_Lifeportal.txt
7: hsa_mature_T.fa
6: HighDL_Chow_8623
5: HighDL_Chow_1x_4154

Name the workflow "Count miRNAs", or something similar.

Select only the six steps from **Clip** to **Cut**. Click **Create Workflow** on the top of the page.

Next, click **Workflow** on the top menu of the LifePortal screen. Then, click on the "Count miRNAs" workflow and select **Edit**.

UiO University of Oslo Analyze Data **Workflow** Shared Data Visualization Apply for a project Help User Using 41.4 GB

Lifeportal

Your workflows

Create new workflow Upload or import workflow

Name	# of Steps
Count miRNAs	6
M Run	5
W Share or Publish history 'imported: LifePortal demo (cheating...)'	28
U Download or Export	3
L Copy	8
W View history 'LifePortal demo 1'	3

Workflows shared with you by others

No workflows have been shared with you.

Other options

Configure your workflow menu

You will now see a graphical representation of the workflow:

The screenshot displays the Lifeportal workflow editor interface. The top navigation bar includes 'UiO University of Oslo', 'Analyze Data', 'Workflow', 'Shared Data', 'Visualization', 'Apply for a project', 'Help', 'User', and 'Using 41.4 GB'. The main header shows 'Lifeportal' and 'Workflow Canvas | Count miRNAs'. On the left, a 'Tools' sidebar lists various bioinformatics tools. The central 'Workflow Canvas' shows a sequence of steps: 'Clip' (circled in red), 'FastQC:Read QC', 'Bowtie2', 'BAM-to-SAM', 'Count', and 'Cut'. The 'Clip' step is connected to 'FastQC:Read QC', which connects to 'Bowtie2', and so on. The right-hand 'Details' panel for the 'Clip' tool is open, showing options like 'Library to clip', 'Minimum sequence length', 'Source', and 'Enter custom sequence'. The 'Enter custom sequence' dropdown is set to 'Set at runtime'.

Drag and rearrange the boxes in a way that looks nice. Try to understand what the boxes and lines represent, as well as the input and output arrows.

A workflow is not connected to specific datasets, but may be reused for any dataset. As the 3' adapter sequence will differ between datasets, the choice of adapter sequence should be up to the user. To configure this:

Click the box for the **Clip** job.

On the right-hand side, click the arrow above **Enter custom clipping sequence** and select **Set at runtime**.

To save the workflow, click the **cogwheel** icon  and select **Save**.

Now, we want to run the workflow on the last HCHF dataset:

Click the **cogwheel** icon  and select **Run**.

Uio University of Oslo Analyze Data Workflow Shared Data Visualization Apply for a project Help User Using 41.4 GB

Lifeportal

Tools
search tools
Upload file
Tests
Bioportal Phylogeny Tools
R
Autodock
Uio tools
Gaussian
NCBI BLAST+
Text Manipulation
Filter and Sort
Join, Subtract and Group
Convert Formats
Extract Features
DeFuse
Bismark
Fetch Sequences
Operate on Genomic Intervals
Metagenomic analyses
Statistics
BEDTools
Graph/Display Data
FASTA manipulation
NGS: GATK2
NGS: Picard (beta)
NGS: Trim Galore
NGS: Assembly
Multiple Alignments
NGS: QC and manipulation
NGS: Mapping
NGS: RNA Analysis
NGS: SAM Tools
NGS: CATK Tools (beta)

Running workflow "Count miRNAs"
Expand All Collapse
History

Projects/Accounts
You have logged as user : sveinung.gundersen@medisin.uio.no. Please select the project (account) where you will be running your workflow.
nrr

Step 1: Clip (version 0.0.13)
Library to clip
3: HighLDL_HCHF_8623
Minimum sequence length (after clipping, sequences shorter than this length will be discarded)
17
Source
Enter custom sequence
Enter custom clipping sequence
ATCTCGTATGCCGTCTTCTGCTTGT
enter non-zero value to keep the adapter sequence and x bases that follow it
0
Discard sequences with unknown (N) bases
Yes
Output options
Output only clipped sequences (i.e. sequences which contained the adapter)

JOB PARAMETERS

Number of tasks
Enter the number of tasks for the current job. (Enter 1, if not sure.)

Number of tasks per node
Enter the number of tasks per node (max 16). (Leave 1, if not sure.)

Walltime (job duration)
Select how long the job shall be running :

Memory per CPU
Select how much memory you need for the job (the allowed format is digit+gb/mb, e.g. 16GB or 3000Mb):

Select a project.

Under the **Clip** job, select the next HCHF input file ("3: HighLDL_HCHF_8623"). Enter the custom clipping sequence: ATCTCGTATGCCGTCTTCTGCTTGT

Under the **Bowtie2** step, make sure that the correct reference genome is selected.

Click **Run workflow** at the bottom of the page. Now, it is time to take a break, as it will take some minutes to run the workflow.

7. Create a full workflow that also joins the counts for all three input datasets

At this point we have the counts for the different experiments in different files. To simplify subsequent analysis, it is better to have only two files. One with the counts for HCHF and one with the counts for Chow. The first step is to join the HCHF count files.

You will add a new step for this, but in this exercise, you will add the step via the workflow editor. Also, you should now be quite skillful in using Galaxy, so you will get a little less step-by-step input at this stage.

First, extract a new workflow "Count 3 miRNAs datasets - joint output" (or similar), based upon all six steps, for all three datasets, in total eighteen steps.

In the workflow editor, align the boxes in an orderly fashion, so that the output files (the **Cut** boxes) are beside (perhaps above?) each other. The **Auto Re-layout** option in the **cogwheel** menu might be of help.

Click the tool **Column join** under the category **Join, Subtract and Group**. Move the box near to the **Cut** boxes. Click **Add new Additional Input** to create a third input arrow for the new tool.

Connect the output of the three **Cut** boxes to the inputs of the **Column Join** box.

Type **c1** in the first option box (as the **hinge** column). Type **c1,c2** in **include these columns**. Select **Yes** for Fill empty columns. Select **Single fill value**, and type **0** (the number) as the **Fill value**.

Save the workflow and exit the editor.

Before running the workflow on the Chow data, you must carry out the **Column join** step on the HCHF data manually. Just be sure to exit the workflow editor, click **LifePortal** on top of the screen to see the tool meny. Then, select the tool, and choose the correct input as described above, and click **execute**.

Rename the resulting history item by dicking the **pencil** icon, type "HCHF_count_proper" in the **Name** box, and click **Save**.

Click the **eye** icon on the history item to look at the results.

Now, run the complete workflow using the 3 Chow datasets. Make sure that the correct adapter sequence is selected for all three files. Change the name of the last result to "Chow_count_proper"

Note that many of the jobs now run in parallel, meaning that the complete workflow should now take about the same time as when only one dataset was used.

8. Run a differential expression analysis with EdgeR

We now have count files for both the HCHF samples and the Chow samples. The differential expression analysis will be done with the help of an r-script.

Click on the **R programming language for statistical computing** under **R**.

Input files should be **One R script and data files**.

The **R-script** is called **r_script_Lifeportal.txt**.

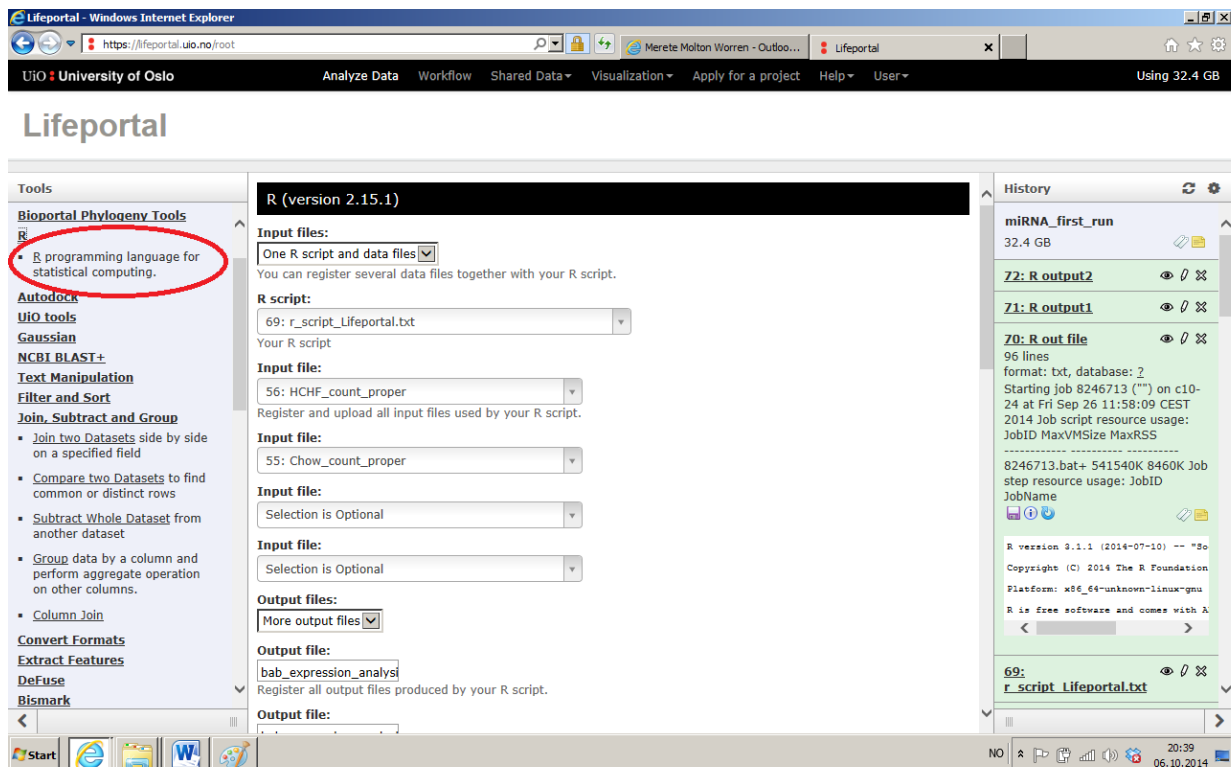
The first input file is the HCHF count file, the next is the Chow count file, the others are left blank.

Output files should be set to **More output files**.

The script produces two output files, and you need to write in their names.

The first is called: **bab_expression_analysis_top_results.txt**

The second is called: **bab_expression_analysis_top_counts.txt**



9. Make a Galaxy page

A Galaxy page is a public page explaining the reasoning and details behind your analysis, together with a presentation of your main results. Most importantly for *reproducibility* is the inclusion of previously run histories, with all related *metadata*, in addition to workflows that other researchers can use to run the analysis on their own datasets.

To create a page select **Saved Pages** from the **User** menu at the top of the LifePortal screen. Then click **Add new page**.

Give the page a meaningful **title** and description (**page annotation**), and click **submit**.

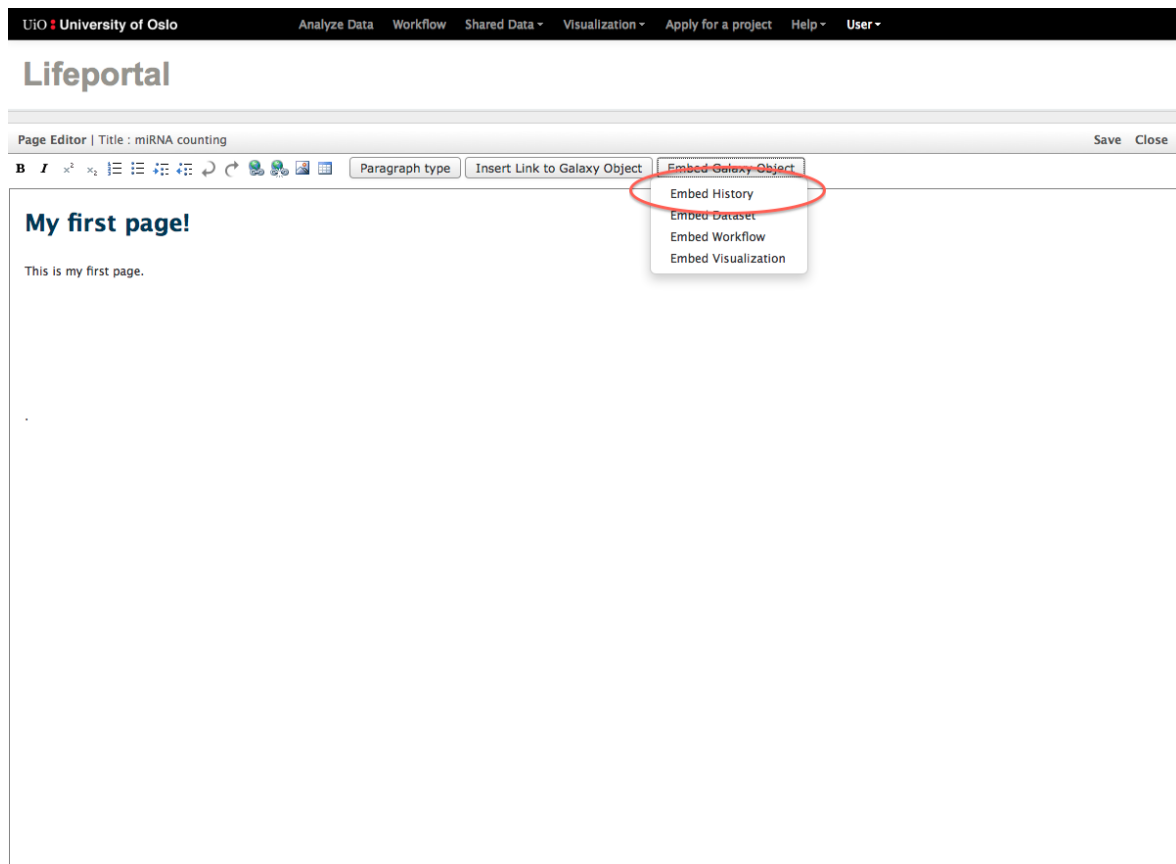
In the list of pages, click the name of the new page and select **Edit contents**.

Create a nice looking page with information about the analysis you have just carried out. Use the **Paragraph type** menu to create headings as needed.

Embed your history and the last workflow using the choices under **Embed Galaxy object***.

*Note: The implementation of Galaxy Pages contains a bug complicating the creation of pages (resulting in content not being displayed in the output). In our experience, some simple measures avoid the issue: First, when beginning a page, add several blank lines and end the document with a period: "." After having

embedding histories or other elements, make sure to manually jump to the next pre-entered line (by clicking) instead of creating a new line with the newline (return) key on the keyboard.



When you are finished, click **Save** and **Close**.

Click the page title and select **Edit** to view the finished results.

Lastly, share your page by clicking the page title selecting **Share** or **Publish**.

Click **Make Page Accessible via Link** to get a **URL** (internet address) you can share with people on email or put in your article.

Click **Share with a user** to share the page to the teachers of this course (meretewo@uio.no and sveinugu@uio.no).