# EMC Galaxy Training - 1: Introduction to Galaxy

This practical aims to familiarize you with the Galaxy user interface. It will teach you how to perform basic tasks such as uploading data, running tools, working with histories, creating workflows, and sharing data.

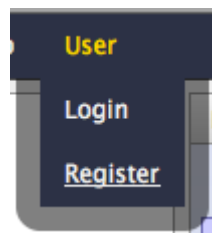**Preparations**

**1. Open Galaxy**
Please open a web browser and navigate to your assigned Galaxy server:

galaxy-training1.trait-ctmm.cloudlet.sara.nl
galaxy-training2.trait-ctmm.cloudlet.sara.nl
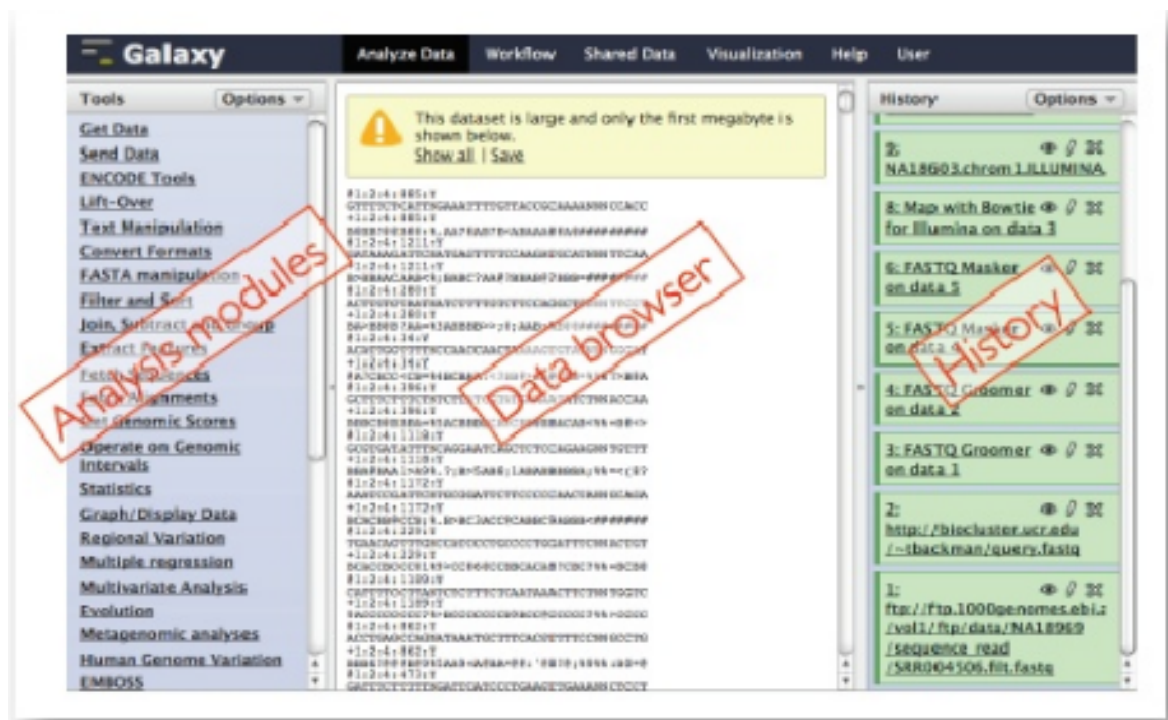galaxy-training3.trait-ctmm.cloudlet.sara.nl
.. etc

It is recommended that you use Firefox or Chrome for this tutorial (not IE)

**2. Register for an account**
In the top menu bar, go to *User* and then choose *Register*. After registration, click on *Analyze data* in the top menu to return to the main screen.

The main screen consists of three parts, on the left is the list of available tools, on the right side is your *history pane*, showing the analysis you have performed so far, and in the middle you view your tools and data.



Follow the steps on the following pages adapted from the *Galaxy 101* course made by the Galaxy Team and one of many tutorials available from the Galaxy Teaching Resources wiki (https://wiki.galaxyproject.org/Teach/Resources).

# Galaxy 101: The first thing you should try

In this very simple example we will introduce you to bare basics of Galaxy:

• Getting data from UCSC
• Performing simple data manipulation
• Understanding Galaxy's History system
• Creating and editing workflows
• Applying workflows to your data

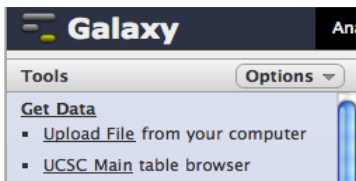You can watch a step-by-step explanation of this entire tutorial here.

## What are we trying to do?

Suppose you get the following question: "Mom (or Dad) ... Which coding exon has the highest number of single nucleotide polymorphisms on chromosome 22?".  You think to yourself "Wow! This is a simple question ... I know exactly where the data is (at UCSC) but how do I actually compute this?" The truth is, there is really no straightforward way of answering this question in a time frame comparable to the attention span of a 7-year-old. Well ... actually there is and it is called Galaxy. So let's try it...
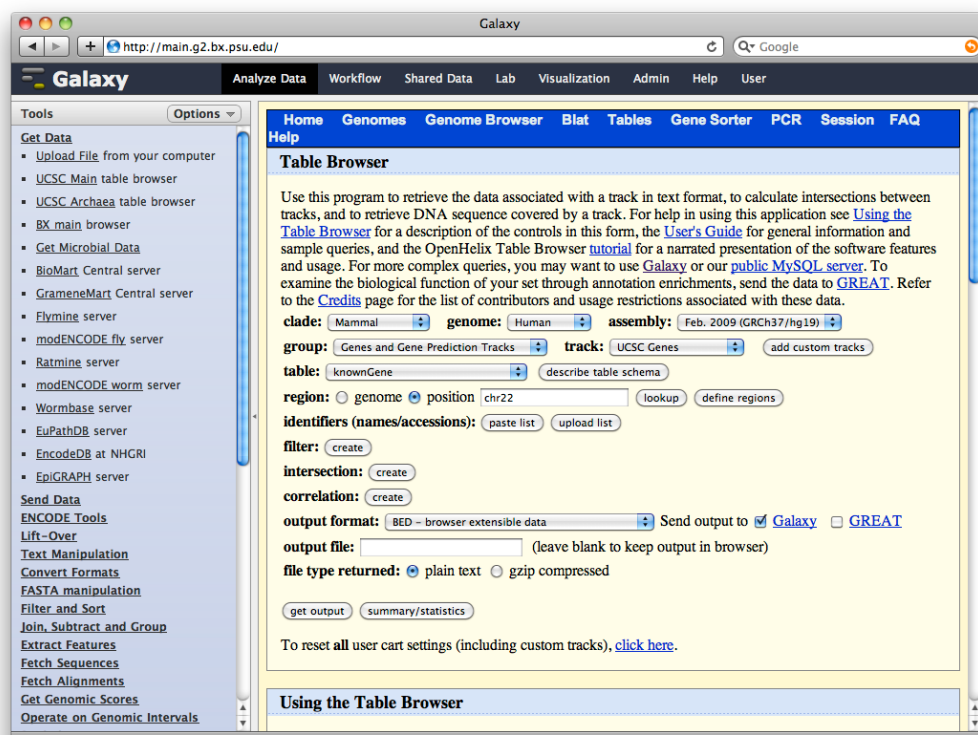
# 1. Getting data from UCSC

## 1.0. Getting coding exons

First thing we will do is to obtain data from UCSC by clicking "Get Data -> UCSC Main":



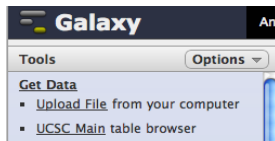You will see Galaxy's middle pane change to looks like this:



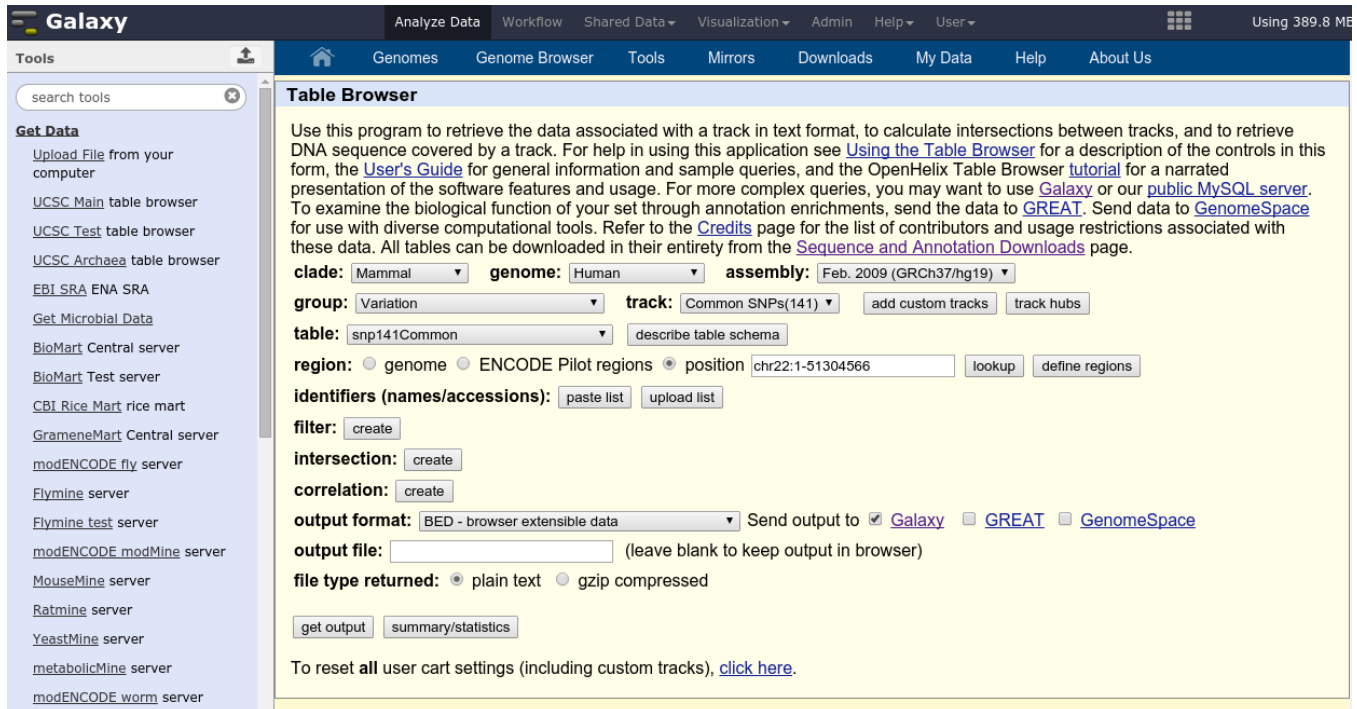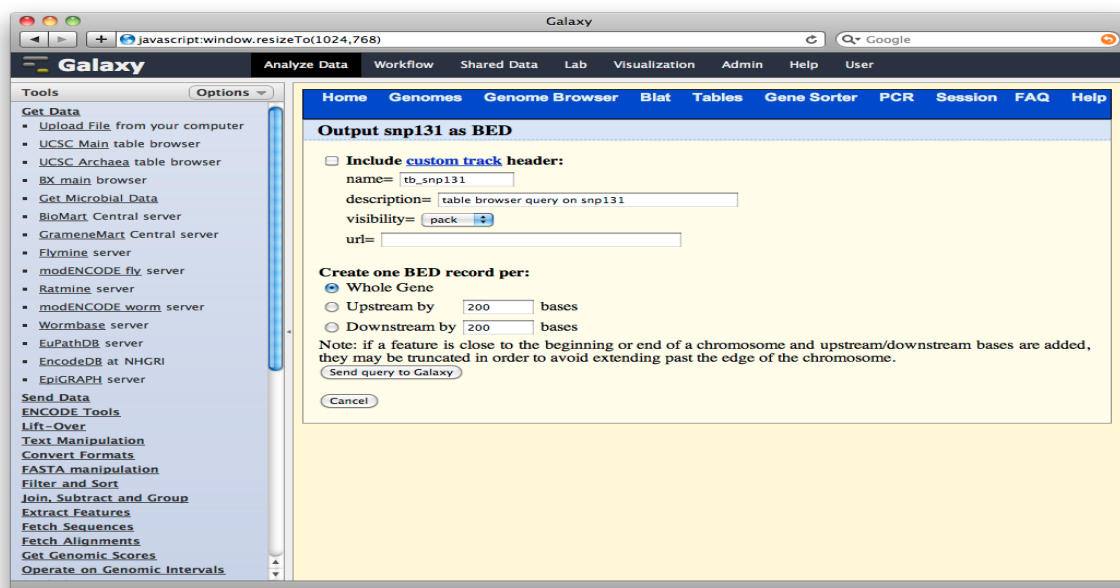Make sure that your settings are exactly the

same as shown on the screen (in particular, **position** should be set to "chr22", **output format** should be set to "BED - browser extensible data", and "Galaxy" should be checked by **Send output to** option). Click **get output** and you will see the next screen:



here make sure **Create one BED record per** is set to "Coding Exons" and click **Send Query to Galaxy**. After this you will see your first History Item in Galaxy's right pane. It will go through gray (preparing) and yellow (running) states to become green:



To view the contents of the file, click on the eye icon.

# 1.1. Getting SNPs

Now is the time to obtain SNP data. This is done almost exactly the same way. First thing we will do is to again click on "Get Data -> UCSC Main":



but now change **group** to "Variation" so that the whole page looks like this:



click **get output** and you should see this:



where you need to make sure that **Whole Gene** is selected ("Whole Gene" here really means "Whole Feature") and click **Send Query to Galaxy**. You will get your second item in the history:

Now we will rename the two history items to "Exons" and "SNPs" by clicking on the Pencil icon adjacent to each item. Also we will rename history to "Galaxy 101" (or whatever you want) by clicking on "Unnamed history" so everything looks like this:



**NOTE:** If the import from UCSC takes too long, the files can also be found in a shared data library in Galaxy. Browse to "Shared Data → Data Libraries", select the data library named "Training Day 1: Introduction to Galaxy", select the files named "SNPs" and "Exons", choose "Import to current history" and click "go". Click on "Analyze Data" on the top menu bar to return to your analysis.



**Data Library "Training Day 1: Introduction to Galaxy"**    [Add datasets] [Add folder] [Library Actions]

✓ 2 datasets imported into 1 history: Unnamed history

| Name | Message | Data type | Date uploaded | File size |
|------|---------|-----------|---------------|-----------|
| ☑ Exons ▾ | None | bed | Wed Nov 12 11:45:53 2014 (UTC) | 779.0 KB |
| ☑ SNPs ▾ | None | bed | Wed Nov 12 11:45:54 2014 (UTC) | 6.8 MB |

For selected datasets: [Import to current history ▾] [Go]

ⓘ TIP: You can download individual library datasets by selecting "Download this dataset" from the context menu (triangle) next to each dataset's name.

ⓘ TIP: Several compression options are available for downloading multiple library datasets simultaneously:

- gzip: Recommended for fast network connections
- bzip2: Recommended for slower network connections (smaller size but takes longer to compress)
- zip: Not recommended but is provided as an option for those who cannot open the above formats
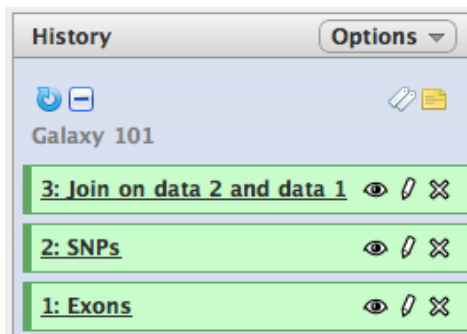
# 2. Finding Exons with the highest number of SNPs

## 2.0. Joining exons with SNPs

Let's remind ourselves that our objective was to find which exon contains the most SNPs. This first step in answering this question will be joining exons with SNPs (a fancy word for printing exons and SNPs that overlap side by side). This is done using the tool  "Join, Subtract and Group -> Join the intervals of two datasets side-by-side**"



**Note** that if you scroll down on this page, you will find an explanation about the tool.

make sure your exons are first and SNPs are second and click **Execute**. You will get the third history item:



which will contain the following data:

…
```
chr22 16277747 16277885 uc002zlj.1_cds_4_0_chr22_16277748_r 0 - chr22 16277851 16277852 rs200742649 0 +
chr22 16287253 16287390 uc002zlj.1_cds_8_0_chr22_16287254_r 0 - chr22 16287338 16287339 rs199952431 0 +
chr22 16287253 16287390 uc002zlj.1_cds_8_0_chr22_16287254_r 0 - chr22 16287345 16287346 rs200013113 0 +
chr22 16287253 16287390 uc002zlj.1_cds_8_0_chr22_16287254_r 0 - chr22 16287371 16287372 rs201840700 0 +
chr22 16448823 16449804 uc011agd.2_cds_0_0_chr22_16448824_r 0 - chr22 16449404 16449405 rs201956705 0 +
```
…

Let's take a look at this dataset. The first six columns correspond to exons. The last six correspond to SNPs. You can see that exon with ID `uc002zlj.1_cds_8_0_chr22_16287254_r` contains three SNPs with IDs  rs199952431, rs200073113, and rs201840700.

# 2.1. Counting the number of SNPs per exon

Above we've seen that exon `uc002z1j.1_cds_8_0_chr22_16287254_r` is repeated three times in the above dataset. Thus we can easily compute the number of SNPs per exon by simply counting the number of repetitions of name for each exon. This can be easily done with the "Join, Subtract, and Group -> Group" tool:



choose column 4 by selecting "c4" in **Group by column**. Then click on **Add new Operation** and select a counting operation on column c4. Set it to ignore lines beginning with #. Make sure the interface looks exactly as shown below:



click **Execute**. Your history will look like this:

if you look at the above image you will see that the result of grouping (dataset #4) contains two columns. This first contains the exon name while the second shows the number of times this name has been repeated in dataset #3.

## 2.3. Sorting exons by SNP count

To see which exon has the highest number of SNPs we can simply sort the dataset #4 on the second column in descending order. This is done with "Filter and Sort -> Sort":



This will generate the fifth history item:

and you can now see that the highest number of SNPs per exon is 30.

## 2.4. Selecting top five

Now let's select top five exons with the highest number of SNPs. For this we will use "Text Manipulation -> Select First" tool:

Clicking **Execute** will produce the sixth history item that will contain just five lines:

## 2.5. Recovering exon info and displaying data in genome browsers

Now we know that in this dataset the five top exons contain between 16 and 30 SNPs. But what else can we learn about these? To know more we need to get back the positional information (coordinates) of these exons. This information was lost at the grouping step and now all we have is just two columns. To get coordinates back we will match the names of exons in dataset #6 (column 1) against names of the exons in the original dataset #1 (column 4). This can be done with "Join, Subtract and Group -> Compare two Datasets" tool (note the settings of the tool in the middle pane):



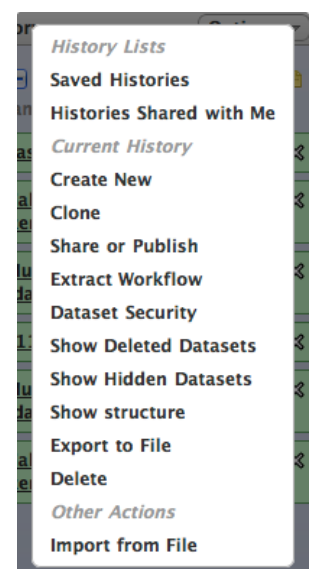this adds the seventh dataset to the history:

The best way to learn about these exons is to look at their genomic surrounding. There is really no better way to do this than using genome browsers. For example, clicking on **"display at UCSC main"** will show this (to see your regions look at "User Supplied Track" (track near the top)):

Enter coordinates: **chr22:32,102,914-32,118,375** to get the same view as in the screenshot below. This centers around the second exon in your top-5 list.



# 3. Understanding histories

In Galaxy your analyses live in histories such as your current one. Histories can be very large, you can have as many histories as you want, and all history behaviour is controlled by the **Options** button on the top of the History pane (gear symbol):

If you create a new history, your current history does not disappear. If you would like to list all of your histories just choose Saved Histories and you will see a list of all your histories in the center pane:



# 4. Converting histories into workflows

One of the history options listed above is very special. It allows you to easily convert existing histories into analysis workflows. Why would you want to create a workflows out of a history? To redo the analysis again with minimal clicking.

## 4.0. Extracting workfklow

Lets take a look at the history again:
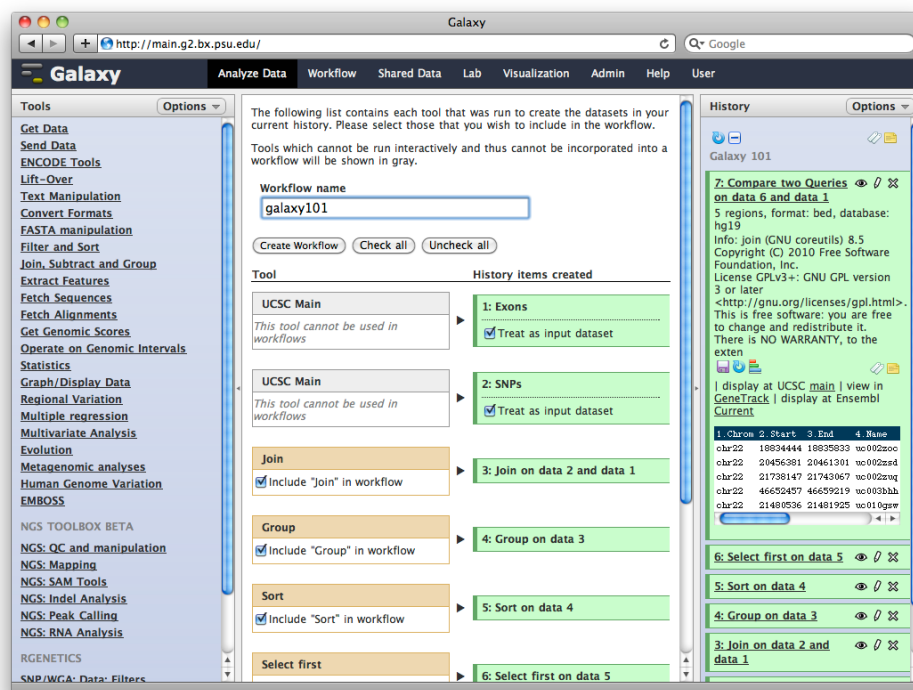
You can see that this history contains all steps of our analysis. So by building this history we have actually built a complete record of our analysis with Galaxy preserving all parameter settings applied at every step. Wouldn't it be nice to just convert this history into a workflow that we'll be able to execute again and again? This can be done by clicking on **Options** (gear) button and selecting **Extract Workflow** option:
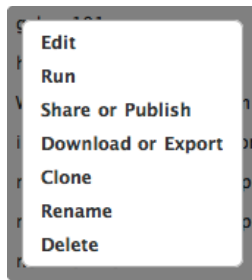


The center pane will change as shown below and you will be able to choose which steps to include/exclude and how to name the newly created workflow. In this case I named it "galaxy101":
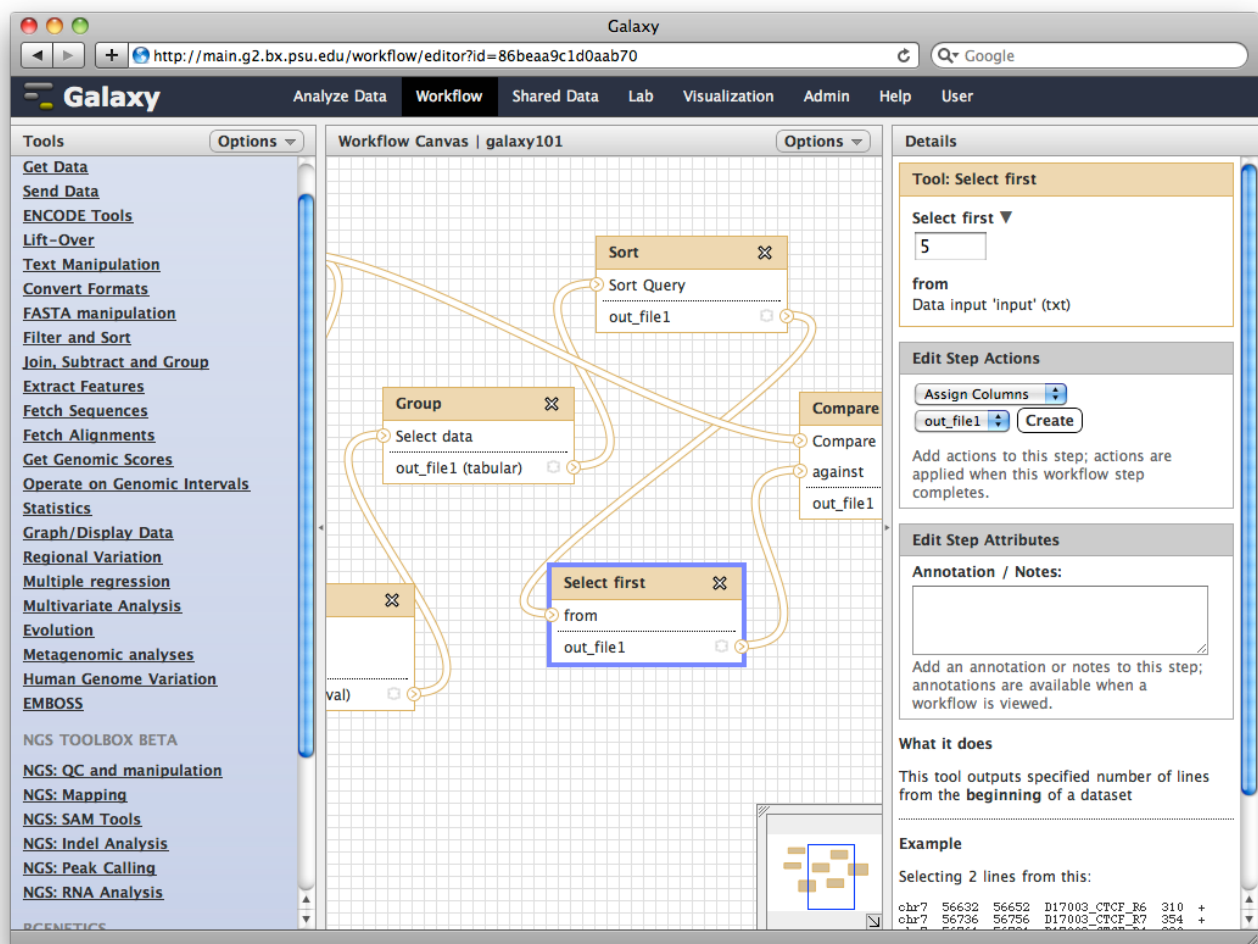


once you click **Create Workflow** you will get the following message: "Workflow 'galaxy101' created from current history". But where did it go? Click on **Workflow** link at the top of Galaxy interface and you will a list of all workflows with "galaxy101" listed at the top:

# 4.1. Opening workflow editor

If you click on a triangle adjacent to the workflow's name you will see the following dialogue:

Edit
Run
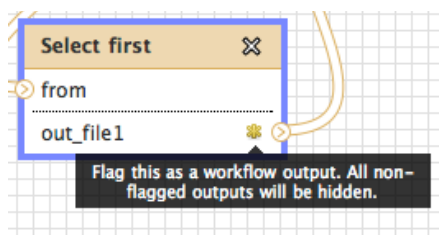Share or Publish
Download or Export
Clone
Rename
Delete

Click **Edit** and the workflow editor will launch. It will allow you to examine and change settings of this workflow as shown below. Note that the box corresponding to the "*Select First*" tool is selected (highlighted with the blue border) and you can see parameters of this tool on the right pane. This is how you can view and change parameters of all tools involved in the workflow.
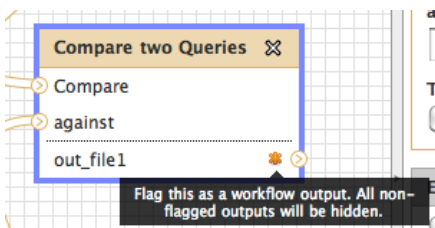


# 4.2. Hiding intermediate steps

Among multiple things you can do with workflows I will just mention one. When workflow is executed one is usually interested in the final product and not in the intermediate steps. These steps can be hidden by mousing over a small asterisk in the lower right corner of every tool box:
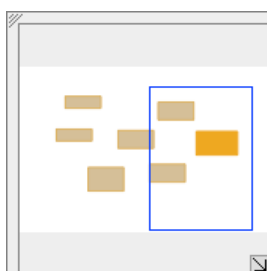
Yet there is a catch. In a newly created workflow all steps are hidden by default and default behavior of Galaxy is that if all steps of a given workflow are hidden, then nothing gets hidden in the history. This may be counterintuitive, but this is done to decrease the amount of clicking if you do want to hide some steps. So in our case if we want to hide all intermediate steps with the exception of the last one we will click that asterisk in last step of the workflow:
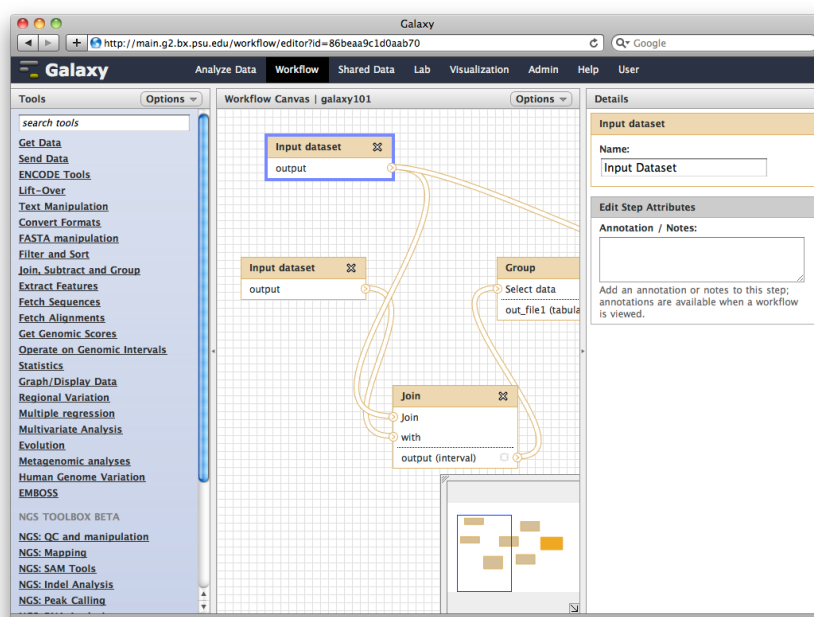


Once you do this the representation of the workflow in the bottom right corner of the editor will change with the last step becoming orange. This means that this is the only step, which will generate a dataset visible in the history:



## 4.3. Renaming inputs

Right now both inputs to the workflow look exactly the same. This is a problem as will be very confusing which input should be exons and which should be SNPs:

One the image above you will see that the top input dataset (the one with the blue border) connects to the *Join* tool first, so it must correspond to the exon data. If you click on this box (in the image above it is already clicked on because it is outlined with the blue border) you will be able to rename the dataset in the right pane:



Then click on the second input dataset and rename it "Features" (this would make this workflow a bit more generic, which will be useful later in this tutorial):



## 4.4. Renaming outputs

Finally let's rename the workflow's output. For this click on the last dataset ("Compare two Queries") and in the **Edit Step Actions** dialogue box select "Rename Dataset"

Click **Create:**



and call it something like "top 5 exons":



## 4.5. Save! It is important...

Now let's save the changes we've made by clicking **Options** (gear symbol) and selecting **Save:**



# 5. Run workflow on different data

Now that we have a workflow, let's use it on some different data.  For example, let's find exons with the highest number of repetitive elements.

### 5.0. Create a new history

Before we start let's create a new history by clicking **Options** and selecting **Create New**.

Now let's get the chr22 exons from the data library again, as well as a list of repetitive elements (which were also obtained from UCSC table browser)

# 5.3. Start the Workflow

At this point you will have two items in your history - one with exons and one with repeats. First, click on the **Workflow** link at the top of Galaxy interface, mouse over "galaxy101" (or whatever you named your workflow), and click on the arrow:



choose **Run.**

Center pane will change to allow you launching the workflow. Select appropriate datasets for Repeats and Exon inputs as shown below, scroll down, and click Run workflow.



Once the workflow has started you will initially be able to see all its steps.



(you may need to click the refresh button at the top of your history if the steps do not show up)

When it is finished, you will see this:



Note that because all intermediate steps of the workflow were hidden, once it is finished you will only see the final dataset #7

If we want to view the intermediate files, we can view the hidden datasets by selecting **"Include Hidden datasets"** from the history options menu.

# 5.5 Share your work

Often you may want to share a workflow you created, or an analysis history, with other users.

You can share a history by clicking on the gear symbol, and selecting **Share or Publish**.

You can do 3 things here:

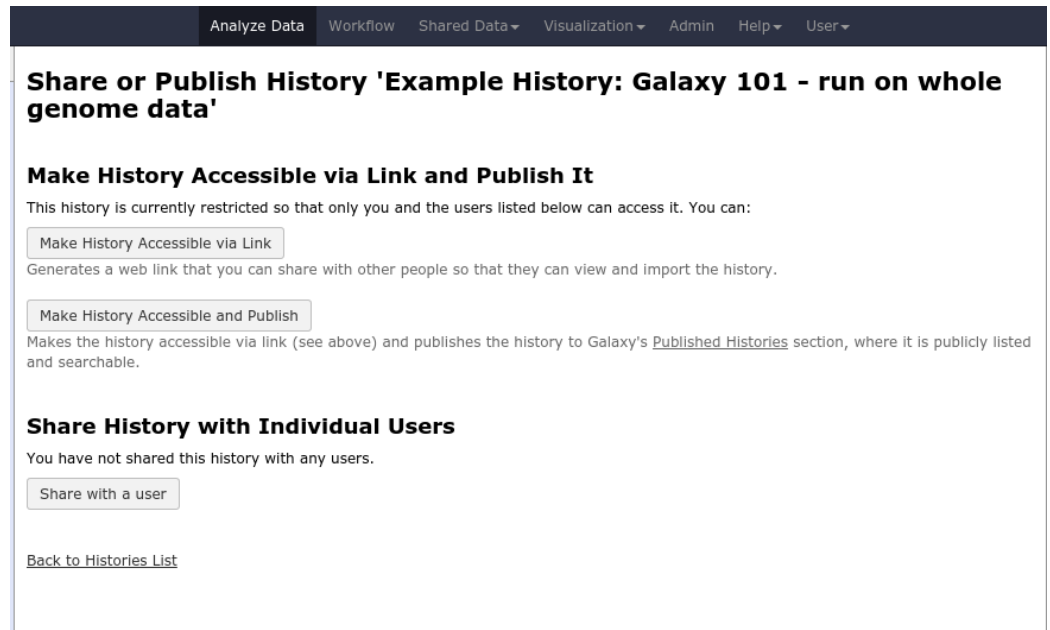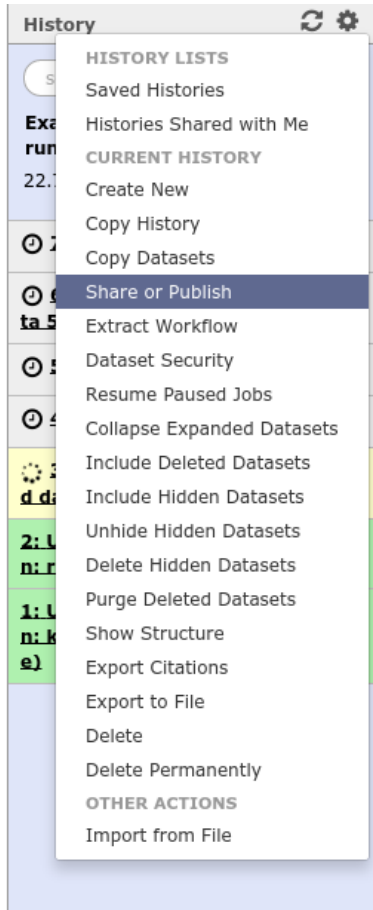> **- Make accessible via Link**
> This generates a link that you can give out to others. Anybody with this link will be able to view your history (even without a Galaxy account)
>
> **- Publish History**
> This will not only create a link, but will also publish your history. (i.e. Your history will now appear under Shared Data → Published Histories)
>
> **- Share with Individual Users**
> This will share the history with specific users on the Galaxy instance. Enter their email address (which they used to register their account in Galaxy)

- Share one of your histories with your neighbour, or publish it.

- See if you can do the same with your workflow!

# THE END!