# Some Phenotype Association Tools in Galaxy: Looking for Disease SNPs in a Full Genome

**Belinda M. Giardine,[1] Cathy Riemer,[1] Richard Burhans,[1] Aakrosh Ratan,[1] and Webb Miller[1]**

[1]Pennsylvania State University, University Park, Pennsylvania

## ABSTRACT

This unit focuses on some of the tools available on the public Galaxy server that are useful for exploring possible associations between human genetic variants and phenotypes. We trace step-by-step through an example illustrating several methods for examining a single full-coverage genome to look for single-nucleotide polymorphisms (SNPs) that are either known to be associated with disease or suspected to have impact for other reasons. It makes use of public genomic data, tools designed specifically for working with variants, and also some general tools for text manipulation and operations on genomic coordinates. *Curr. Protoc. Bioinform.* 39:15.2.1-15.2.27. © 2012 by John Wiley & Sons, Inc.

Keywords: disease • SNP • genome variation • coding • non-coding • gene-based analysis • Web application

## INTRODUCTION

Galaxy (*http://galaxyproject.org*; Giardine et al., 2005; Blankenberg et al., 2010; Goecks et al., 2010) is a software framework that provides Web-based tools for bioinformatics, including tasks useful in the analysis of human variation. The developers maintain a public server at Penn State, and the software is also freely available for local installation and customization. Galaxy is highly extensible, so as new tools become available (not necessarily written specifically for Galaxy), they can be added to increase the power and flexibility of the system.

This unit focuses on some of the tools available on the public Galaxy server that are useful for exploring possible associations between human genetic variants and phenotypes. We trace step-by-step through an example illustrating several methods for examining a single full-coverage genome to look for single-nucleotide polymorphisms (SNPs) that are either known to be associated with disease or suspected to have impact for other reasons. It makes use of public genomic data, tools designed specifically for working with variants, and also some general tools for text manipulation and operations on genomic coordinates. A version of this tutorial is also available in HTML format (see Internet Resources, below). For a more basic introduction to using Galaxy in general, please see UNIT 10.5, "Using Galaxy to Perform Large-Scale Interactive Data Analyses."

For this example, we will use an artificial dataset consisting of the SNP calls from the Complete Genomics genome GS12880 (Drmanac et al., 2009) with a few known disease variants added. This will provide a realistic background to search for the disease SNPs, but not necessarily a realistic collection of disease SNPs for a single individual to have. We chose an assortment of six SNPs from the PhenCode database (Giardine et al., 2007), representing different genes and different parts of the gene (Table 15.2.1). There are two coding SNPs (heterozygous) and four non-coding (one heterozygous and three homozygous). The four non-coding SNPs are located in a promoter region, a UTR, and two introns.

Understanding
Genome
Variation

**Table 15.2.1** Disease SNPs Planted in the Sample Dataset[a]

| chr11 | 5248153 | 5248154 | G | Intron | HbVar:thalassemia |
|-------|---------|---------|---|--------|-------------------|
| chr11 | 5255743 | 5255744 | C | UTR | HbVar:thalassemia |
| chr11 | 5275879 | 5275880 | C/T | Coding | HbVar:Hb E |
| chr16 | 222915 | 222916 | T/C | Coding | HbVar:Hb Lyon-Bron |
| chrX | 31279779 | 31279780 | T/C | Intron | LMDp:muscular dystrophy |
| chrX | 100641248 | 100641249 | G | Promoter | BTKbase:Agammaglobulinemia |

[a]The position of the SNP is given here using the same coordinate conventions as in BED format, to match what is seen in the results. The last column shows the database the SNP is from and the associated disease.

This example builds a single sequential history, but it is organized into several protocols according to the type of analysis being performed. These illustrate the following skills.

*Basic Protocol 1: Preparing input data*

- Uploading files

- Using Galaxy libraries

- Basic filtering

*Basic Protocol 2: Selecting known coding SNPs predicted to be damaging, then finding their genes and associated pathways*

- PolyPhen-2

- Gene-based analysis

*Basic Protocol 3: Running new predictions of coding SNPs likely to be detrimental*

- SIFT

- Using published workflows

*Basic Protocol 4: Finding SNPs that fall in suspected functional regions*

- Predicted regulatory regions

- ENCODE functional data

- PhyloP conserved positions

**Conventions Used in Figures (Screen Shots) for Tutorial Purposes**

Red arrows or boxes on the screen shots indicate settings or things you will need to do. Green arrows indicate the "go" buttons once the settings or parameters are selected. Blue arrows or boxes point out additional information that you should note, but that do not require any action.

**USING GALAXY TO LOOK FOR DISEASE SNPs IN A FULL GENOME: PREPARING INPUT DATA**

*Overview*

- Start with a full set of SNP calls from a particular individual, in the masterVar format used by Complete Genomics.

- Upload the file to a Galaxy history via URL or FTP.

Phenotype
Association Tools
in Galaxy

**15.2.2**

Supplement 39

Current Protocols in Bioinformatics

- Convert to pgSnp format.

- Obtain a set of SNPs found in healthy individuals from the Galaxy library.

- Subtract the SNPs found in healthy individuals to narrow the search.

### Necessary Resources

*Hardware*

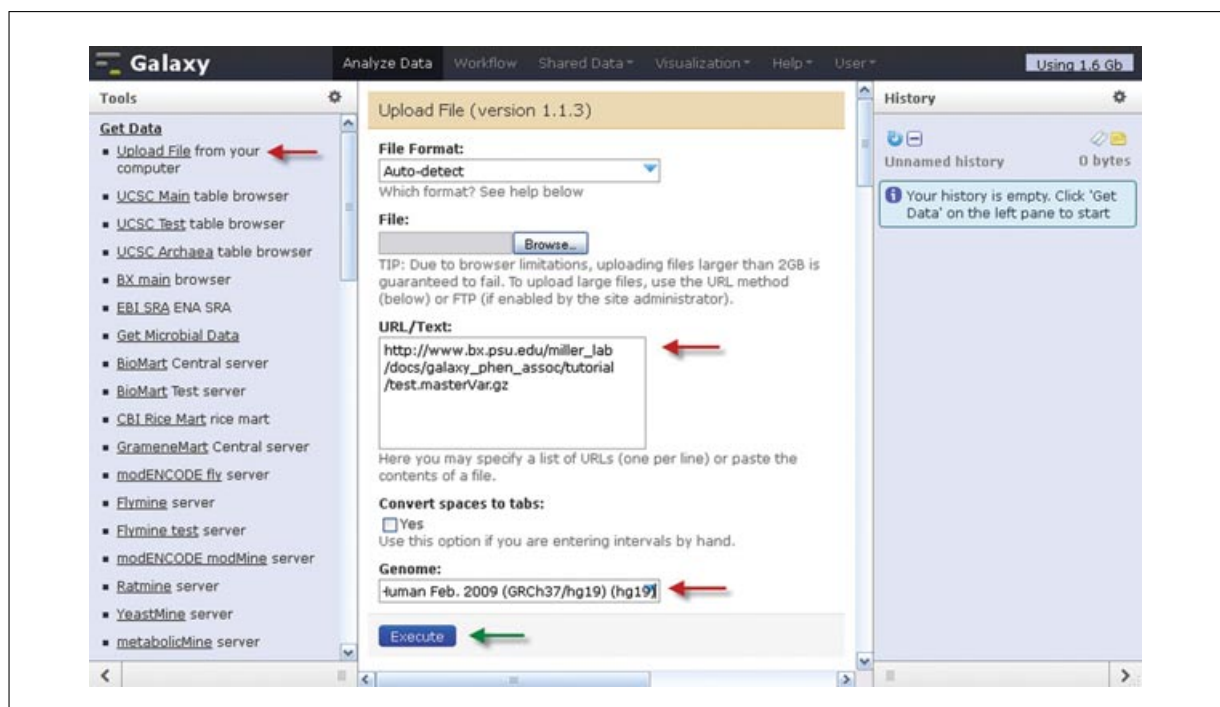An Internet-connected computer

*Software*

Internet browser that supports JavaScript (e.g., most current browsers such as Mozilla Firefox, Safari, Opera, Chrome, or Microsoft Internet Explorer)
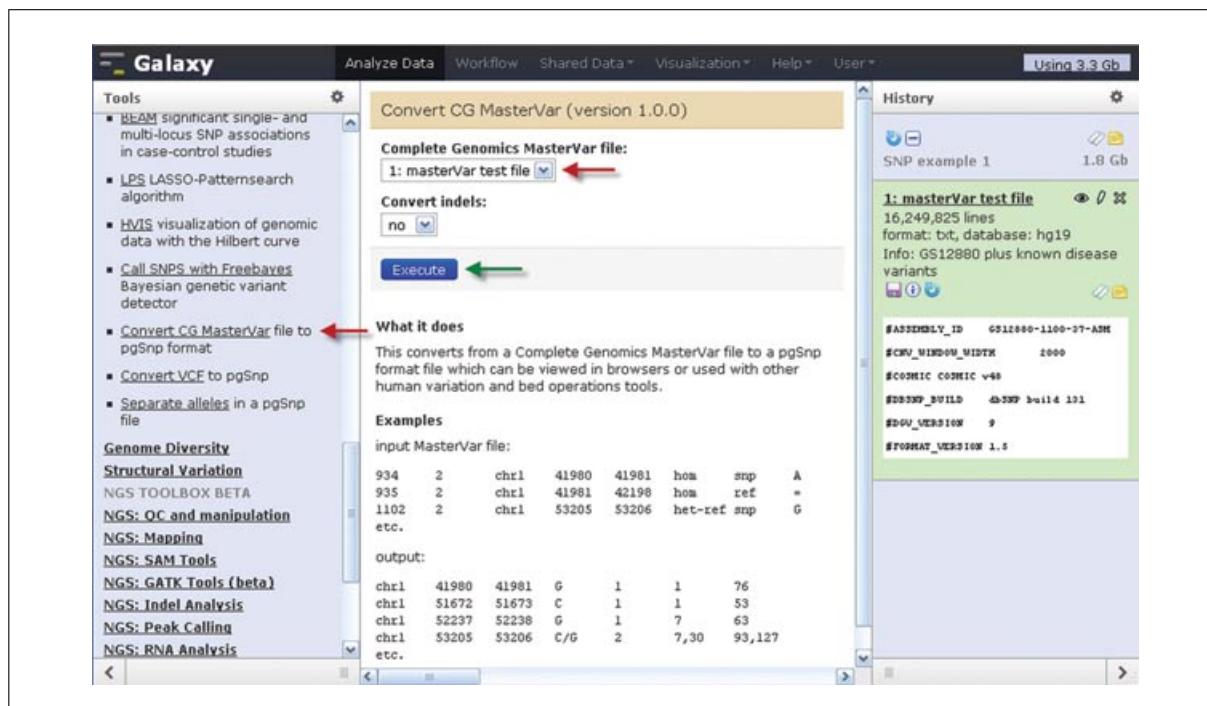
*Files*

Supplementary file: `test.masterVar.gz` (found at *http://www.bx. psu.edu/miller_lab/docs/galaxy_phen_assoc/tutorial/test.masterVar.gz*). See Supplementary File at the end of this unit for an alternate URL to access the supplementary file.

### Uploading a data file

1. There are two ways to upload large files to Galaxy: one is to use FTP, and the other is to place the file in a Web-accessible location and give Galaxy the URL. For this example, we will use the URL method. Open your Web browser, go to the Galaxy portal at *galaxyproject.org*, and choose "Use the free public server." In the tool panel on the left, click on the section called Get Data to open it, then click on the Upload File tool (indicated by a red arrow in Fig. 15.2.1). The input form for the selected tool will appear in the center panel, along with other information about the tool.



**Figure 15.2.1** Uploading a data file. See text for details. For color version of figure go to *http://www.currentprotocols. com/protocol/bi1502*.

**15.2.3**

**Figure 15.2.2** Converting to pgSnp format. See text for details. For color version of figure go to *http://www. currentprotocols.com/protocol/bi1502*.

2. Type or paste the file's URL in the large text box: *http://www.bx.psu.edu/miller_lab/ docs/galaxy_phen_assoc/tutorial/test.masterVar.gz*.

3. Also be sure to choose the corresponding genome assembly in the Genome box. Typing a search term in the select box (e.g., hg19) will limit the list of choices to items containing the specified text. Then, click Execute (green arrow) to upload the dataset into your history.

   *Compressed files are unpacked automatically.*

### Converting to pgSnp format

4. Our uploaded dataset is in the masterVar file format used by Complete Genomics, but that doesn't work with very many of the Galaxy tools. Therefore we will convert it to pgSnp, which is a specialized BED format for SNPs. In the tool panel, open the Phenotype Association section and click on the tool called Convert CG MasterVar (red arrow in Fig. 15.2.2). Then, in the center panel, select the dataset with the masterVar file (red arrow), and click the Execute button (green arrow).
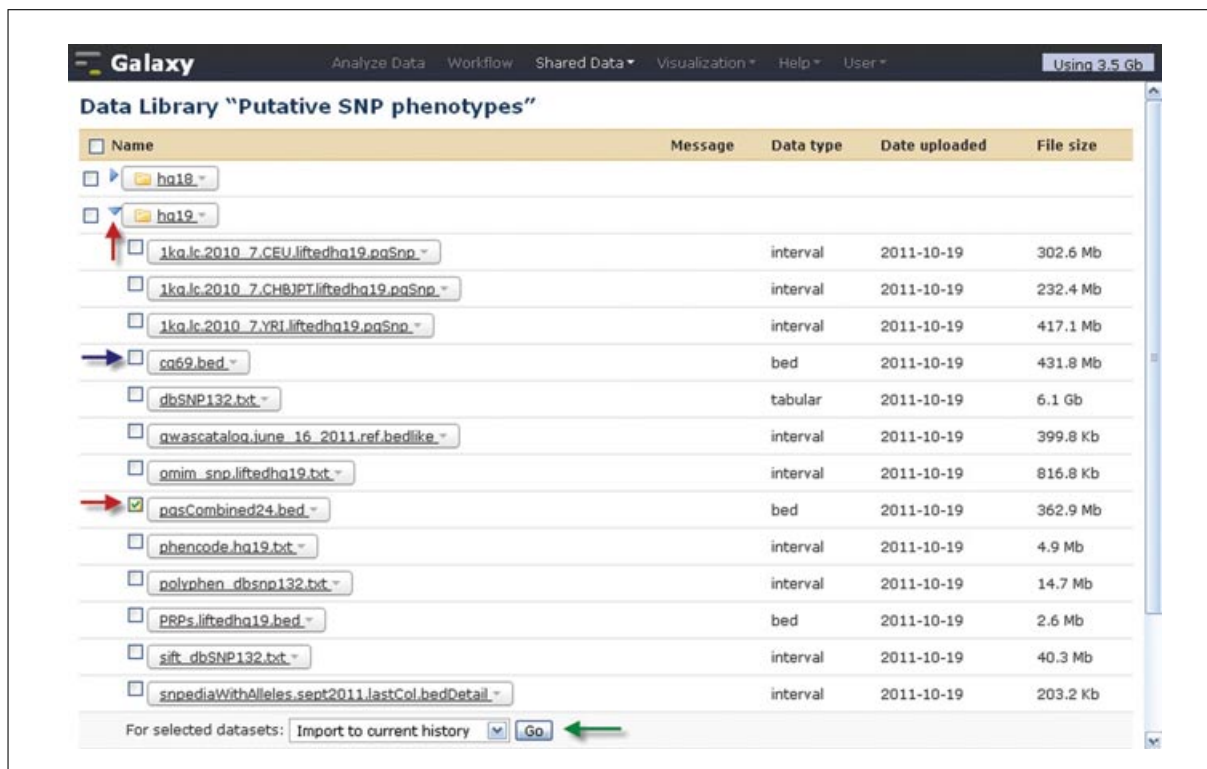
   *It is good practice to rename the major datasets in your history to remind yourself later what each one represents. Here we have renamed the uploaded dataset, and also the history itself (see UNIT 10.5).*

### Galaxy libraries

5. To look for dominantly inherited disease SNPs, it is helpful to remove SNPs found in healthy individuals. To do this we will pull in some data from a Galaxy library. Click on Shared Data in the menu bar and select Data Libraries. Then, look for a library called Putative SNP Phenotypes and click on it.

### Putative SNP Phenotypes library

6. First open the folder with the hg19 datasets by clicking on the blue triangle next to the folder name (red arrow pointing up in Fig. 15.2.3). To import datasets into your history, check the ones you want (red arrow) and then click the Go button (green arrow).

**Figure 15.2.3** Putative SNP Phenotypes library. See text for details. For color version of figure go to *http://www.currentprotocols.com/protocol/bi1502*.

> *In this library, there are currently two sets of full-coverage SNPs from healthy individuals. One has 24 public genomes from a variety of sources, populations, and sequencing technologies (red arrow). The other is a group of 69 from Complete Genomics (blue arrow). In most cases, you will want to use both of these for filtering; however since our example input dataset is based on one of the CG genomes, we will skip that set here.*

### Removing SNPs found in healthy individuals

7. Clicking on the Analyze Data link in the menu bar returns us to the page with the history and tool panels.

8. Now we are ready to remove the "benign" SNPs from our input. In the tool panel, open the section called Operate on Genomic Intervals. Click on Subtract (red arrow in Fig. 15.2.4) to see the options for that tool in the center panel. We want to remove the SNPs that were found in the 24 genomes from our converted set (red arrows). Next, click Execute (green arrow).

   > *The tools in the Operate on Genomic Intervals section are examples of tools that cannot be run on a masterVar file; they only work with "interval" formats, which include pgSnp.*
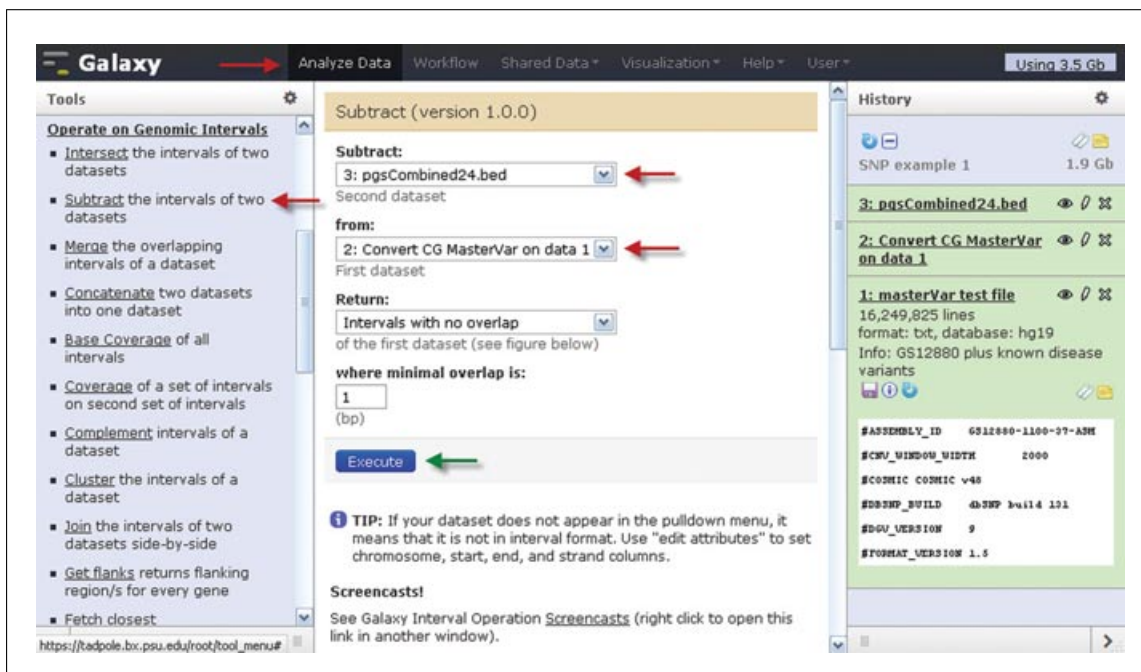
   > *If we had not used a Complete Genomics dataset as our starting point, the next step would be to repeat the subtraction using the result from the first subtraction and the CG dataset from the library (blue arrow in Fig. 15.2.3).*
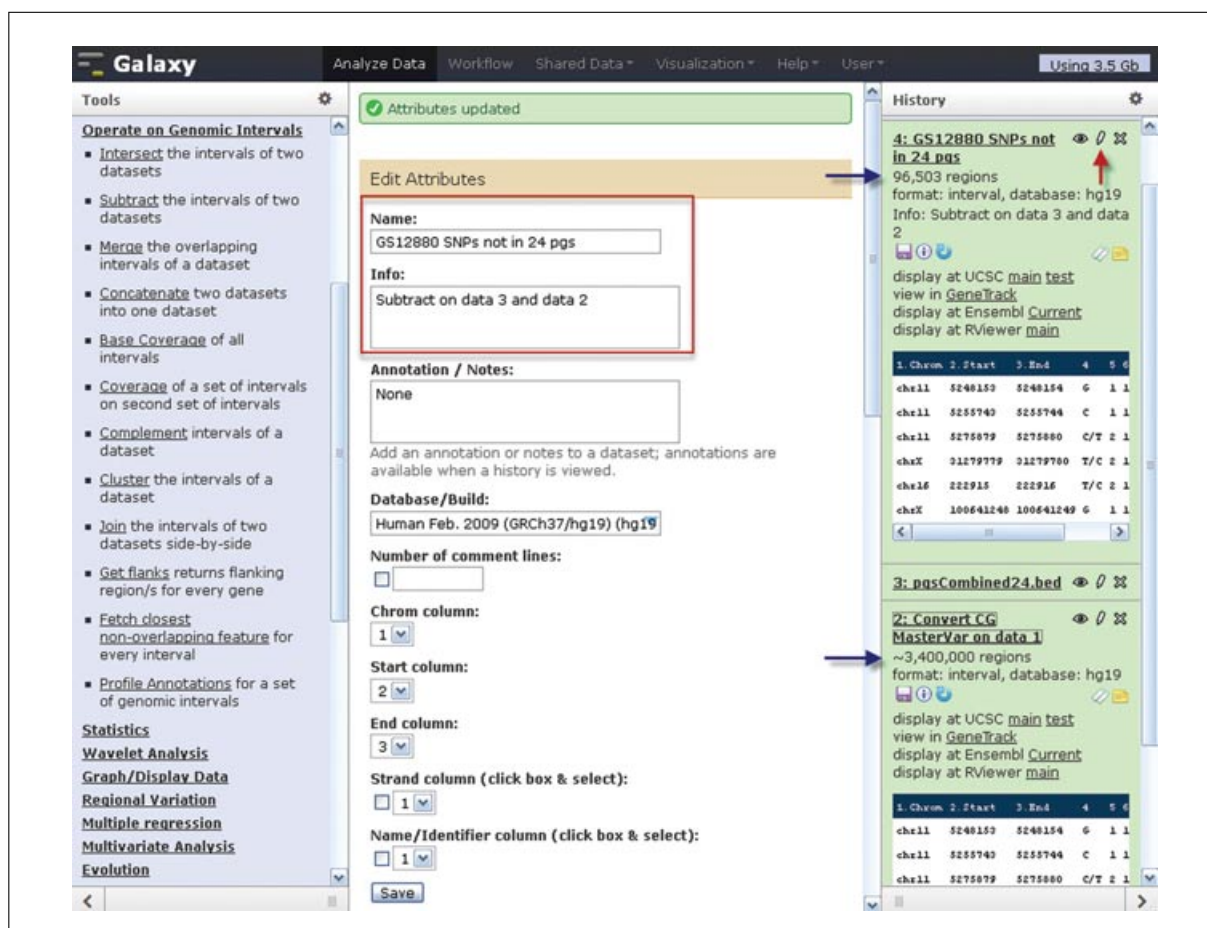
### Completed input dataset

9. Open the results in the history panel by clicking on the dataset name. We see that the number of SNPs is greatly reduced, from around 3.4 million to 96 thousand (blue arrows in Fig. 15.2.5).

   > *We have renamed the results of this step, moving the automatically generated name to the info field (red box). This dataset will be the starting point in the subsequent protocols for this example. The filtering we have done is relevant when looking for both coding and non-coding SNPs.*

**Understanding Genome Variation**

**15.2.5**

**Figure 15.2.4** Removing SNPs found in healthy individuals. See text for details. For color version of figure go to *http://www.currentprotocols.com/protocol/bi1502*.



**Figure 15.2.5** Completed input dataset. See text for details. For color version of figure go to *http://www.currentprotocols.com/protocol/bi1502*.

**15.2.6**

# SELECTING KNOWN CODING SNPs PREDICTED TO BE DAMAGING, THEN FINDING THEIR GENES AND ASSOCIATED PATHWAYS

## *Overview*

● Import a public library file containing pre-computed results from running the PolyPhen-2 program (Adzhubei et al., 2010) on the dbSNP database.

● Join the input dataset with the PolyPhen-2 results row-by-row, based on interval overlap. This adds new columns to the input set, including the UniProt protein accession ID and the predicted effect of the SNP.

● Filter the results to select rows containing the word "damaging."

● Translate the UniProt IDs to HUGO gene symbols by joining with an identifier table imported from UCSC.

● Run the CTD tool to extract curated pathway associations for these genes from the Comparative Toxicogenomics Database.

## *Necessary Resources*

### *Hardware*

An Internet-connected computer

### *Software*

Internet browser that supports JavaScript (e.g., most current browsers such as Mozilla Firefox, Safari, Opera, Chrome, or Microsoft Internet Explorer)

## *PolyPhen-2 predictions*

1. Go to Shared Data > Data Libraries > Putative SNP Phenotypes > hg19 as before. This time look for the dataset `polyphen_dbsnp132.txt`. To learn more about this dataset, click on its name instead of marking the box.

## *Details about the PolyPhen-2 dataset*

2. Most of the first section is information automatically computed about the dataset, including its size, date uploaded, and a peek at the actual data (Fig. 15.2.6). The Message at the top, as well as the second section, Other Information, is provided by the person who added the file to the library. Here we find labels for the data columns, and learn that dbSNP version 131 was used as input to compute the predictions, which are in column 11.

3. To add this dataset to your history from here, click on its name (green arrow) and select "Import this dataset into selected histories." On the next page that appears, your current history should already be selected as the default destination, so just click the "Import library datasets" button. Then, click Analyze Data in the menu bar to return to the page with the history and tool panels.

## *Joining on genomic intervals*

4. We will do a join between the two datasets to get the predictions associated with our SNPs. We do not have a shared identifier to join on, so instead we will join together rows of the two datasets whenever their positions on the genome overlap. The result will have all of the information from both datasets, for any position on the genome where there are data in both datasets.

    In the Operate on Genomic Intervals section of the tool panel, click on Join (Fig. 15.2.7). Then, in the center panel, select the filtered SNPs from Basic Protocol 1 as the first dataset and the PolyPhen-2 predictions as the second. We

**Figure 15.2.6**  Details about the PolyPhen-2 dataset. See text for details. For color version of figure go to *http://www.currentprotocols.com/protocol/bi1502*.
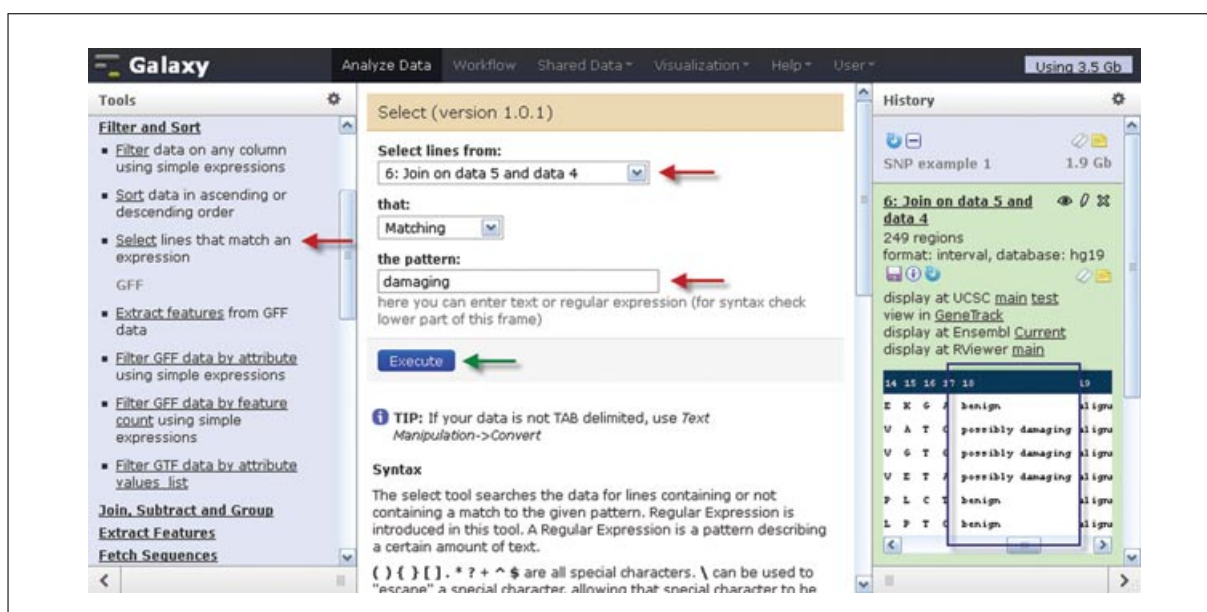
**Figure 15.2.7**  Joining on genomic intervals. See text for details. For color version of figure go to *http://www. currentprotocols.com/protocol/bi1502*.



**Figure 15.2.8**  Selecting damaging results. See text for details. For color version of figure go to *http://www. currentprotocols.com/protocol/bi1502*.

are only interested in SNPs that appear in both sets, so leave the default setting of doing an Inner Join. Then click Execute.

*Selecting damaging results*

5. There are 249 predictions (Fig. 15.2.8), but some of those say that the SNP is benign. In the history panel, a few rows of data are shown; use the scrollbar to scroll over to the predictions (blue box). The ones we are interested in use the word "damaging."

6. Open the Filter and Sort section and choose the Select tool, which selects lines from a dataset that match (or do not match) the given pattern. To run it, choose the dataset with the join results and type in `damaging` for the pattern. The search is case-sensitive, so type it exactly as it appears in the dataset. Then click Execute.

   *The predictions are no longer in column 11 because we requested the columns from the SNP dataset to be listed first in the join results.*

**Understanding Genome Variation**

**15.2.9**

**Figure 15.2.9** PolyPhen-2 results. See text for details. For color version of figure go to *http://www.currentprotocols.com/protocol/bi1502*.



**Figure 15.2.10** Mapping between identifiers. See text for details. For color version of figure go to *http://www.currentprotocols.com/protocol/bi1502*.

### PolyPhen-2 results

7. We find 104 SNPs that are predicted to be damaging by PolyPhen-2 (blue arrow in Fig. 15.2.9), including one of our two coding disease SNPs (not shown). For small datasets like this, you can display the entire contents by clicking on the eye icon in the history panel.

**Figure 15.2.11** Choosing the identifier fields. See text for details. For color version of figure go to *http://www. currentprotocols.com/protocol/bi1502*.

In the blue box you can see that we now have UniProt IDs for many of the genes associated with these variants. We can look for pathways associated with the genes by using the CTD tool; however, that tool requires HUGO/NCBI identifiers (Seal et al., 2011) rather than UniProt IDs for its input, so first we will download a file from the UCSC Table Browser (Karolchik et al., 2004) to map between the identifiers (red arrow in the tool panel).

*The columns do not line up in this view because the data are tab-separated and the values in each column are not all of the same length.*
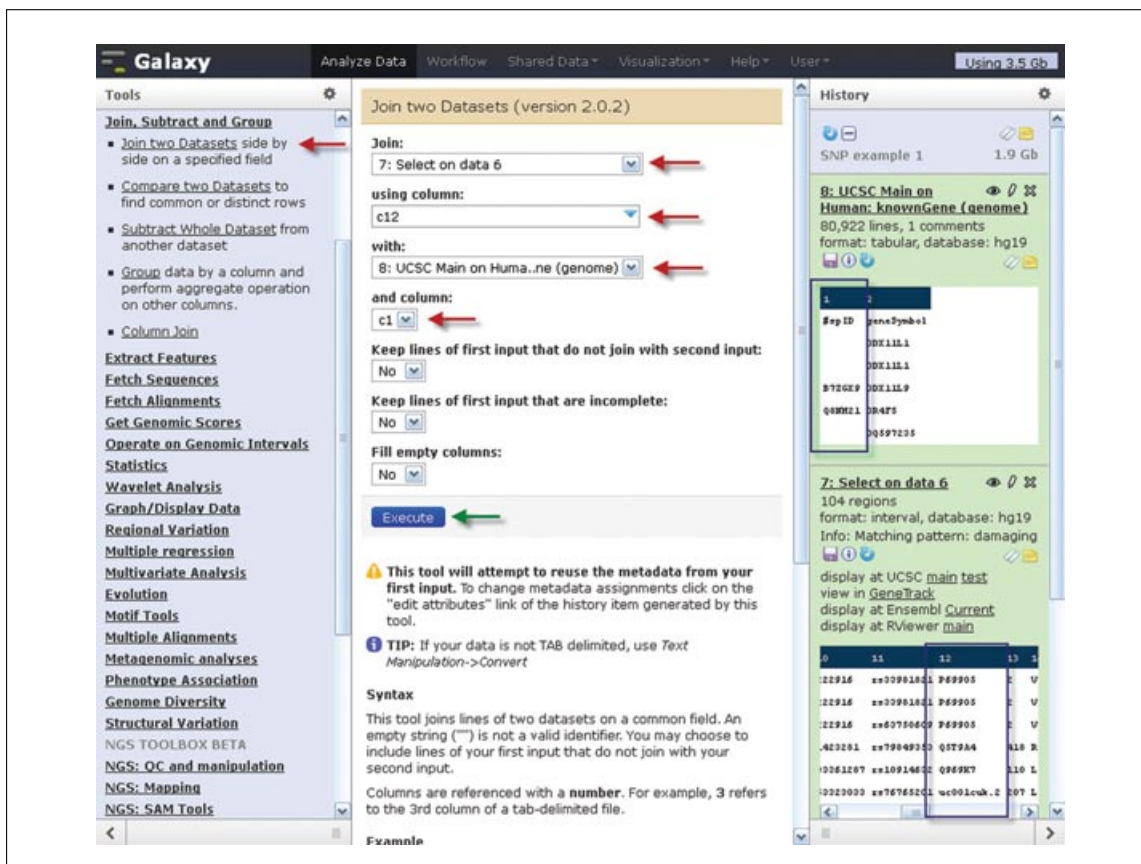
### Mapping between identifiers

8. First make sure the genome and assembly are correct; we want Human build hg19 in order to match our history datasets. We will be using the UCSC Known Genes table because it has the most additional information, including connections between various identifiers, and we want them for the full genome. All of these options are shown in the red box (Fig. 15.2.10). We also want the output format to be "selected fields from primary and related tables," and to have the results sent back to Galaxy (red arrows). Then click the "get output" button (green arrow).

*The box for chromosome coordinates is irrelevant when the radio button for the whole genome is selected.*

*Sometimes having the window divided into three panels makes the center one too narrow to work in easily. In these cases you can click on the large "<" and ">" arrow buttons in the lower corners to hide the tool and/or history panels. When you are done with the*

**Understanding Genome Variation**

**15.2.11**

**Figure 15.2.12** Joining on identifiers. See text for details. For color version of figure go to *http://www. currentprotocols.com/protocol/bi1502*.
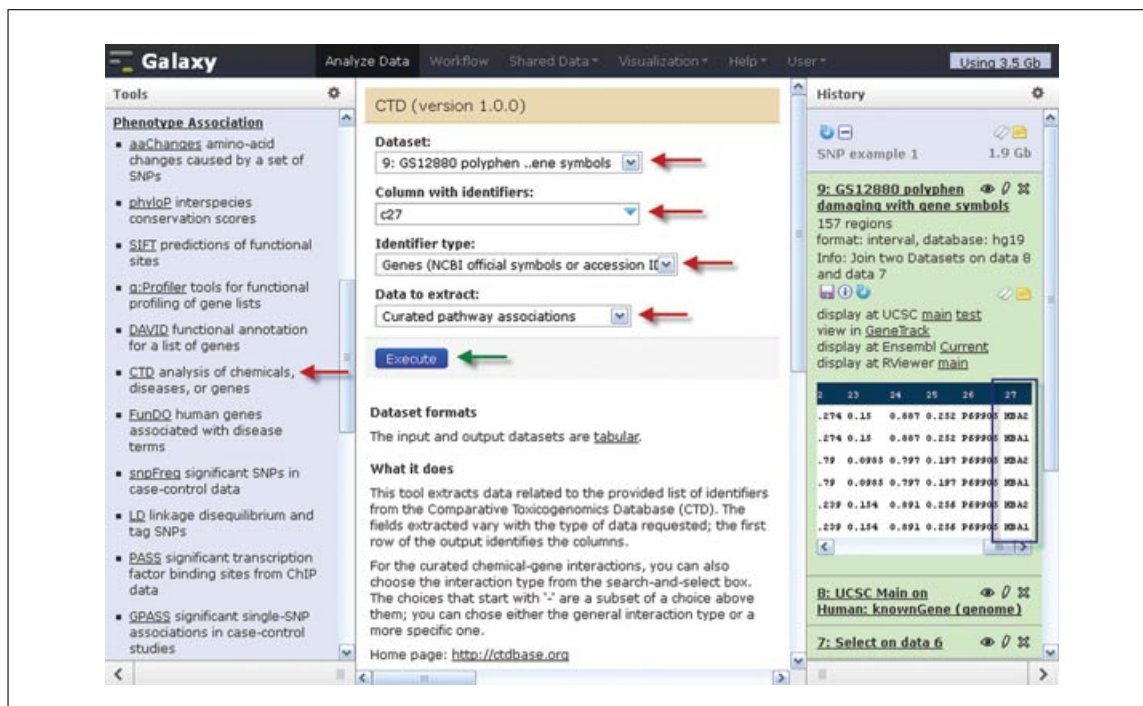
*wider center panel, or need to access something in the other panels, just click the arrows again to reopen them. Here we have closed both panels to more easily see the UCSC Table Browser interface.*
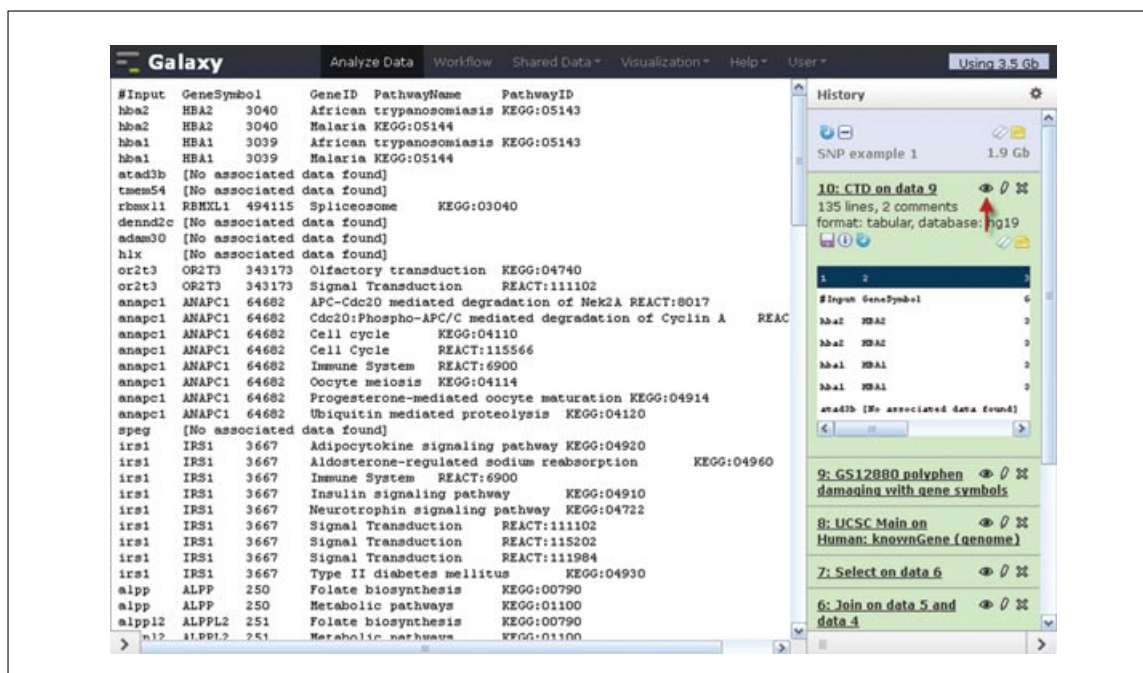
### Choosing the identifier fields

9. Now we select the fields we want (Fig. 15.2.11). The main table does not have what we need, but the cross-reference table below it has both the UniProt ID and the official HUGO gene symbol (red arrows). After marking these boxes, click the "done with selections" button (green arrow). This will bring up one more page with just two buttons; click the one that says "Send query to Galaxy." The Table Browser interface will then close, and the requested data will be sent to your Galaxy history.

### Joining on identifiers

10. Back at Galaxy, we are ready to do a join to get the gene symbols associated with the damaging SNPs. This time we do not have genomic positions in both sets, but they do have a field in common: the UniProt ID. We will do the join by matching up the values in those columns. In the section Join, Subtract, and Group there is a tool called Join Two Datasets (Fig. 15.2.12). The parameters for this tool are the two datasets, the columns to match up, and what to do with unmatched rows. By opening the two datasets in the history you can find the column numbers containing the UniProt IDs (blue boxes). Make the selections shown, putting the SNP dataset first so its fields will be first in the result rows. We are not interested in any of the rows that do not join, so leave the rest of the options set to No. Click Execute.

**Figure 15.2.13** CTD. See text for details. For color version of figure go to *http://www.currentprotocols. com/protocol/bi1502*.



**Figure 15.2.14** CTD results. See text for details. For color version of figure go to *http://www.currentprotocols. com/protocol/bi1502*.

### CTD

11. Now that we have a list of gene symbols for the putatively damaging SNPs, we can look for relationships among these genes. There are several tools useful for this, including CTD (Davis et al., 2009), g:Profiler (Reimand et al., 2007), and DAVID (Huang et al., 2009). For this example we will use the CTD tool to extract pathway information from the Comparative Toxicogenomics Database (Fig. 15.2.13).

**Understanding Genome Variation**

**15.2.13**

Select the join results from the previous step, and column 27 which has the gene symbols (blue box). The identifier type we are using is NCBI gene symbols (which are the same as HUGO), and for the output we want the pathways. Click Execute.

*CTD results*

12. The results are a table of the genes and pathways (Fig. 15.2.14). From here, knowledge of the disease and/or the biology could be used to further filter the SNPs.

BASIC
PROTOCOL 3

**RUNNING NEW PREDICTIONS OF CODING SNPs LIKELY TO BE DETRIMENTAL**

*Overview*

• Augment the input dataset (from Basic Protocol 1) to add a new column containing both the variant and reference alleles, by running the published workflow "Prep pgSnp file to run SIFT."

• Run the SIFT tool, keeping the original allele column as a comment, and requesting the gene name, OMIM disease, etc., in the output.

• Filter the results to select rows containing the word "DAMAGING."

*Necessary Resources*

*Hardware*

An Internet-connected computer

*Software*

Internet browser that supports JavaScript (e.g., most current browsers such as Mozilla Firefox, Safari, Opera, Chrome, or Microsoft Internet Explorer)

*Input for SIFT*

1. PolyPhen-2 is not the only way to predict functional coding changes with Galaxy; the SIFT tool (Kumar et al., 2009) in the Phenotype Association section performs a similar assessment. However, unlike the PolyPhen-2 library dataset which is precomputed only for known SNPs in the dbSNP database, the SIFT tool can make predictions for all possible nucleotide substitutions in the human exome.
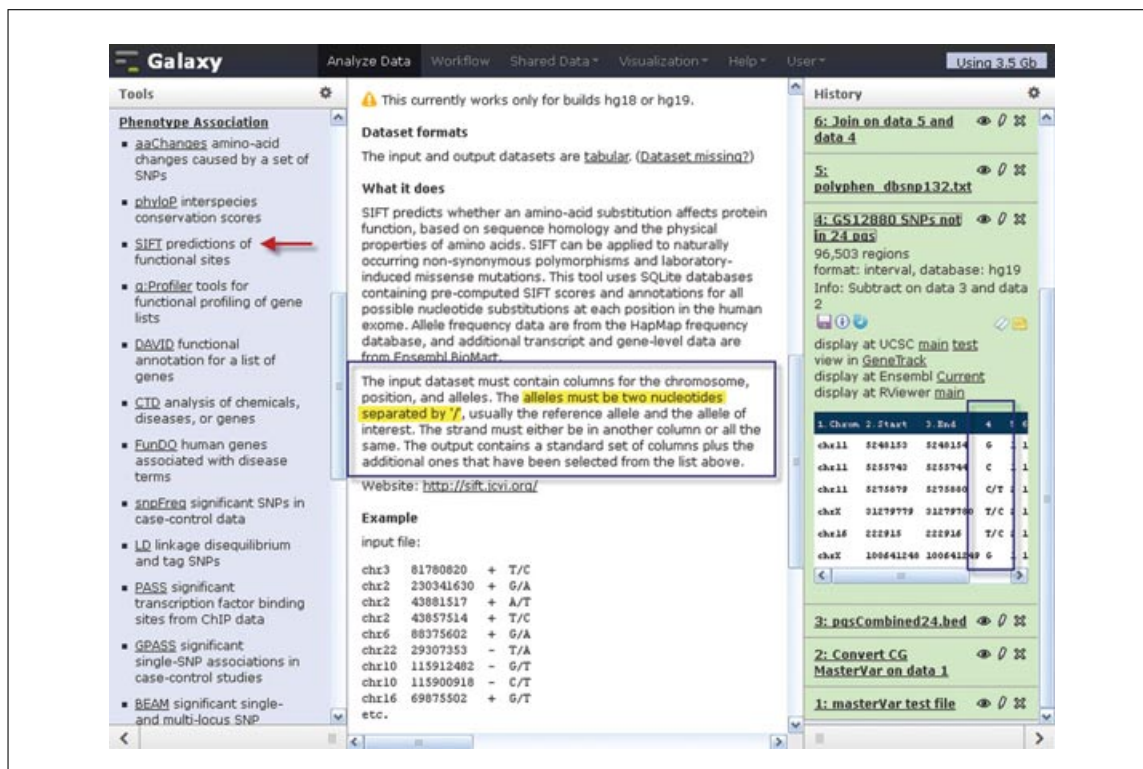
   Each Galaxy tool typically has a description and other helpful information located below its input form, so to learn more about SIFT, open it and scroll down in the center panel as shown in Figure 15.2.15. Here, we find that the format needed (blue box) is a little different than that of our pgSnp dataset; SIFT requires two alleles at each SNP position, but pgSnp format lists only a single nucleotide for homozygous genotypes (blue box in history panel).

2. The SIFT instructions suggest adding the reference nucleotide in such cases; this can be done within Galaxy, but involves quite a few steps. For any complex task that you are likely to repeat, it is helpful to have a *workflow*, which is like an automatic recipe. Happily, there is a published workflow to convert from the pgSnp format to what SIFT needs, so we can do this in one step by using the workflow.
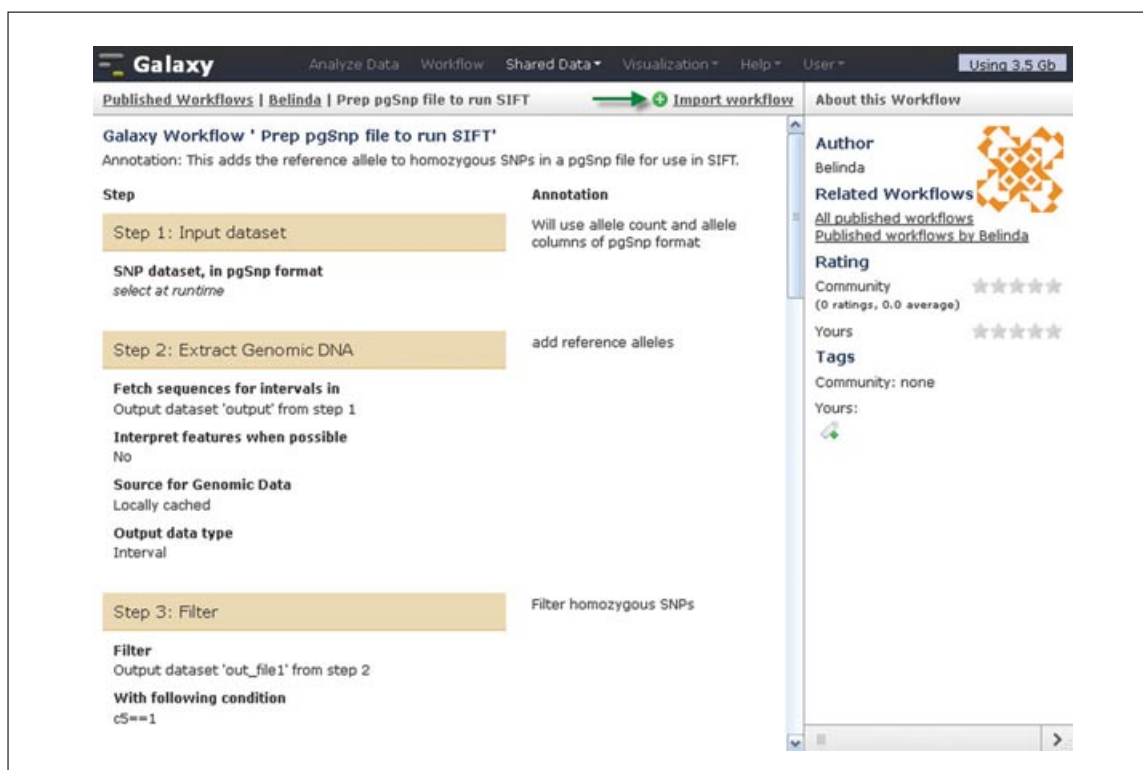
*Finding the workflow*

3. Click on Shared Data in the menu bar and select Published Workflows. In the search box, type SIFT and click on the magnifying glass. Find the workflow "Prep pgSnp file to run SIFT" in the resulting list and click on its name.

Phenotype
Association Tools
in Galaxy

**15.2.14**

Supplement 39

Current Protocols in Bioinformatics

**Figure 15.2.15** Input for SIFT. See text for details. For color version of figure go to *http://www.currentprotocols. com/protocol/bi1502*.



**Figure 15.2.16** Viewing the workflow. See text for details. For color version of figure go to *http://www. currentprotocols.com/protocol/bi1502*.

**Figure 15.2.17** Running the workflow. See text for details. For color version of figure go to *http://www.currentprotocols.com/protocol/bi1502*.

### Viewing the workflow

4. The full workflow is visible, along with annotations (Fig. 15.2.16). There is a description of what the workflow does, and some notes for individual steps. The right panel shows the author, user rating, and tags for this workflow. Click on "Import workflow" (green arrow) to bring this into your personal set of workflows. Then, on the confirmation screen, click "start using this workflow."
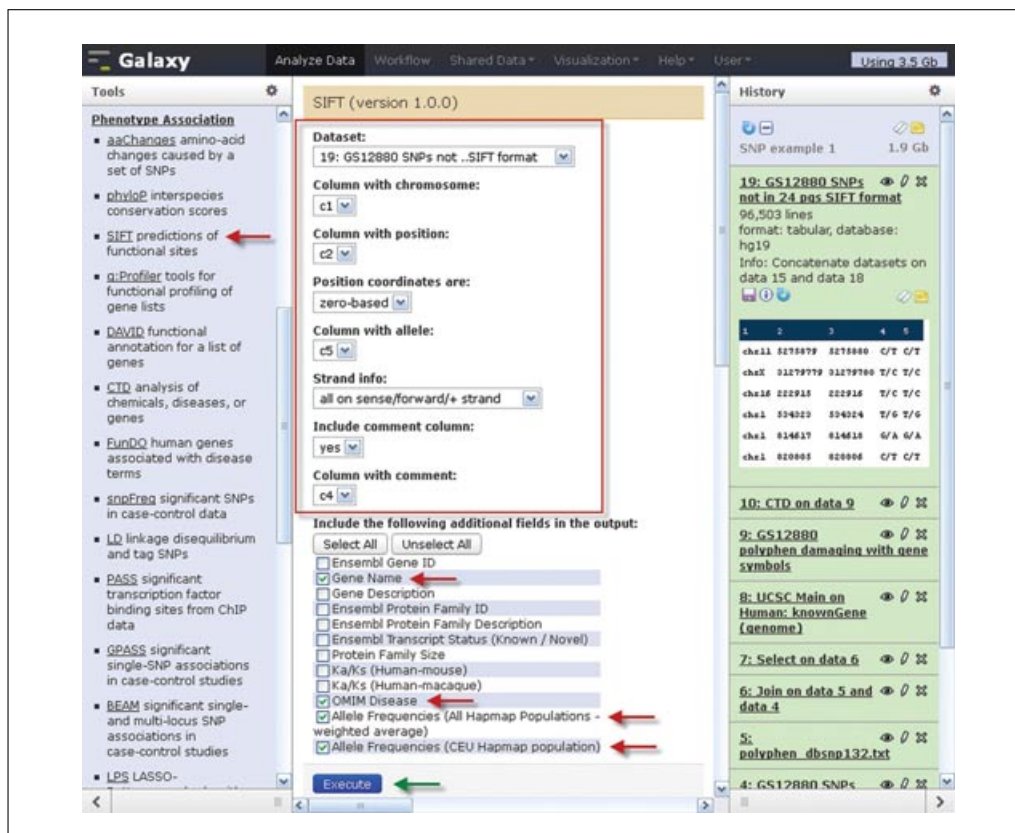
### Your workflows

5. You will be presented with the list of workflows that you have acquired. Click on the name of the one you just imported, and select Run from the menu that pops up.

### Running the workflow

6. The workflow is in the center panel (Fig. 15.2.17); for its input select the same results from Basic Protocol 1 that we used in Basic Protocol 2 for the join with the PolyPhen-2 predictions. Several of the steps have annotations explaining what the workflow will do. The annotation for step 10 (blue arrow) tells us that column 5 of the output will be what we want to submit to SIFT, and that column 4 will have the original allele data. There are ten steps in the workflow, all of which will be run when you click the "Run workflow" button.

### SIFT

7. When the workflow has completed, open the SIFT tool again and fill in the parameters as shown in the red box in Figure 15.2.18. In pgSnp format, the alleles are always on the forward strand, so select "all on sense/forward/+ strand." By choosing "yes" for

**Figure 15.2.18** SIFT. See text for details. For color version of figure go to *http://www.currentprotocols.com/protocol/bi1502*.

"Include comment column" we can keep column 4, which has the original alleles, in the output. We have also selected some of the optional fields to be included in our output (red arrows). Click Execute. This tool may take some time to complete, so please be patient.

*The workflow's intermediate steps are hidden from view as they complete, to make the history less cluttered. If you wish to see all of the steps in your history, click on the gear icon at the top of the history panel. A menu will open up, and one of the choices is Show Hidden Datasets.*
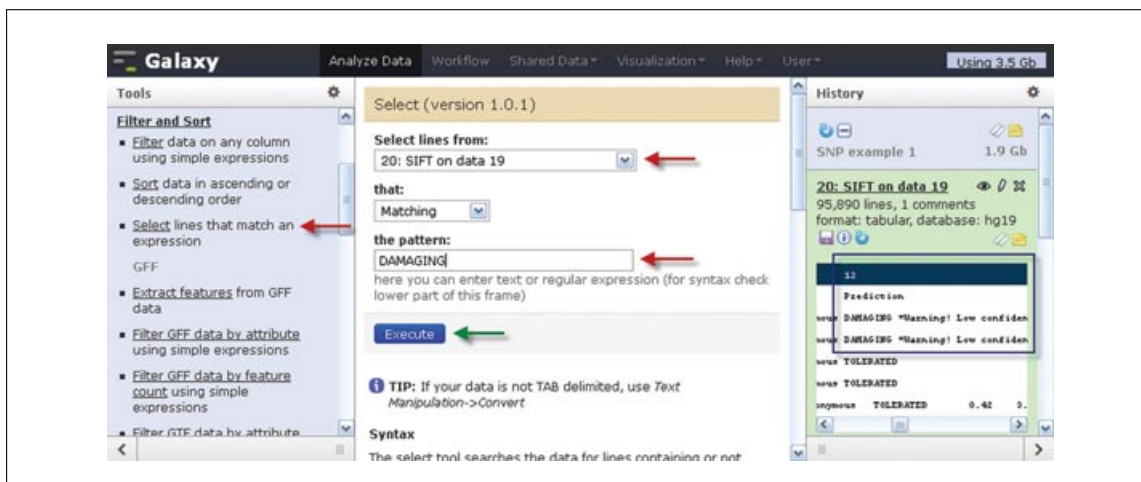
### Selecting damaging SNPs

8. As before with the PolyPhen-2 predictions, we now need to filter the SNPs to select the damaging ones. The SIFT output in the history panel (blue box in Fig. 15.2.19) shows that it uses DAMAGING in upper case. Once the parameters are set, click Execute.
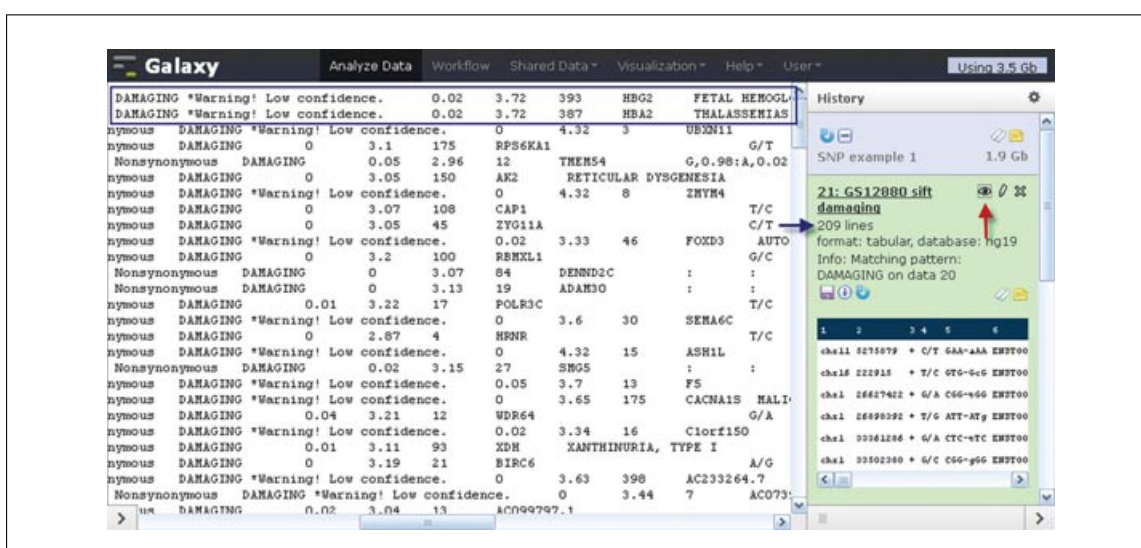
### SIFT results

9. Clicking on the eye icon displays the results in the center panel (Fig. 15.2.20). To see them better, we have temporarily hidden the tool panel by clicking on the arrow button in the bottom left corner. We also renamed this dataset. Within the 209 predicted damaging SNPs, we see that both of our coding disease SNPs are found (blue box). Because we requested the OMIM output field when running SIFT, we have the diseases from OMIM listed here.

*Other tools such as CTD and DAVID could now be run on these genes if desired, but for this example we will move on to hunting for disease SNPs in non-coding regions.*

**Figure 15.2.19** Selecting damaging SNPs. See text for details. For color version of figure go to *http://www.currentprotocols.com/protocol/bi1502*.



**Figure 15.2.20** SIFT results. See text for details. For color version of figure go to *http://www.currentprotocols.com/protocol/bi1502*.

## FINDING SNPs THAT FALL IN SUSPECTED FUNCTIONAL REGIONS

### *Overview*

• Filter the input dataset (from Basic Protocol 1) to keep only rows whose intervals overlap those in a library dataset of predicted regulatory regions.

• In a similar fashion, find rows in the same input dataset that overlap with those in an ENCODE regulatory dataset (DNase clusters) obtained from UCSC.

• Run the PhyloP tool on the same input dataset to add a column of interspecies conservation scores. Then, use the Histogram tool to help choose a suitable score threshold, and filter the SNPs on the score column to keep only those at highly conserved positions.

### *Necessary Resources*

### *Hardware*

An Internet-connected computer

*Software*

Internet browser that supports JavaScript (e.g., most current browsers such as Mozilla Firefox, Safari, Opera, Chrome, or Microsoft Internet Explorer)

*Predicted regulatory regions (PRPs)*

1. First we will import another dataset from the Putative SNP Phenotypes library, this time `PRPs.liftedhg19.bed`, which contains computationally predicted regulatory regions (PRPs). Do not forget to click on Analyze Data after importing the dataset to get back to your history and tool page.
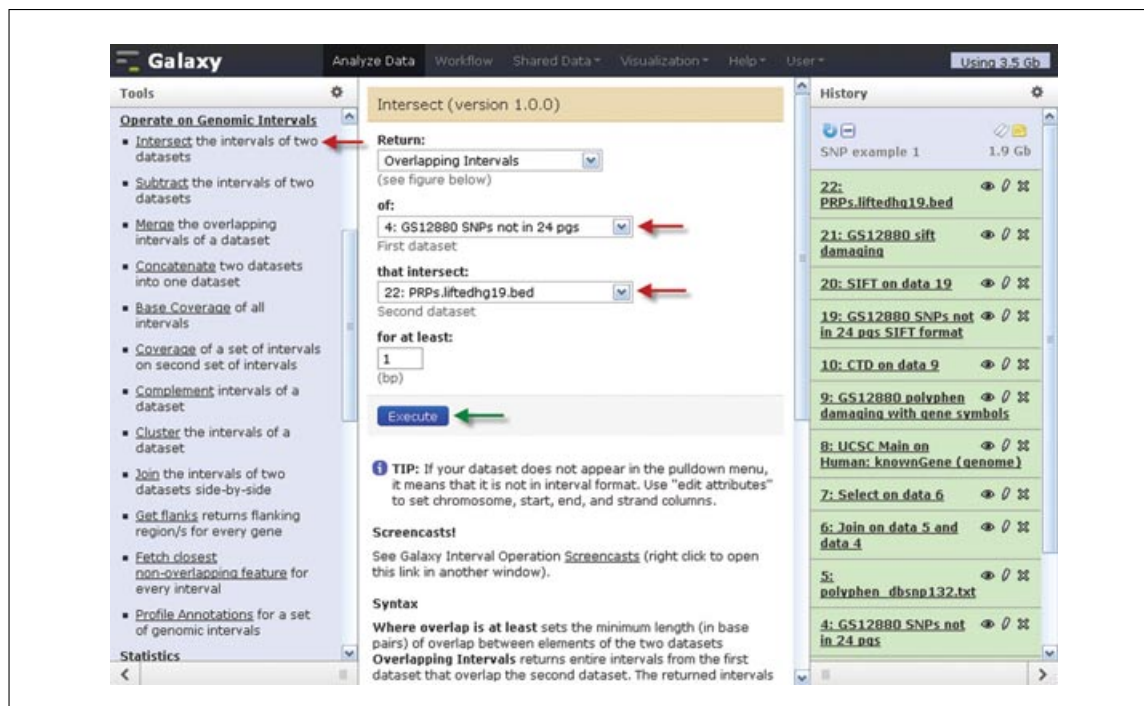
   *This is the intersection of the PreMods (Ferretti et al., 2007) and the regions with high Regulatory Potential scores (Taylor et al. 2006). Both of these sets were computed on earlier genome assemblies and then lifted (i.e., mapped) to the current assembly.*
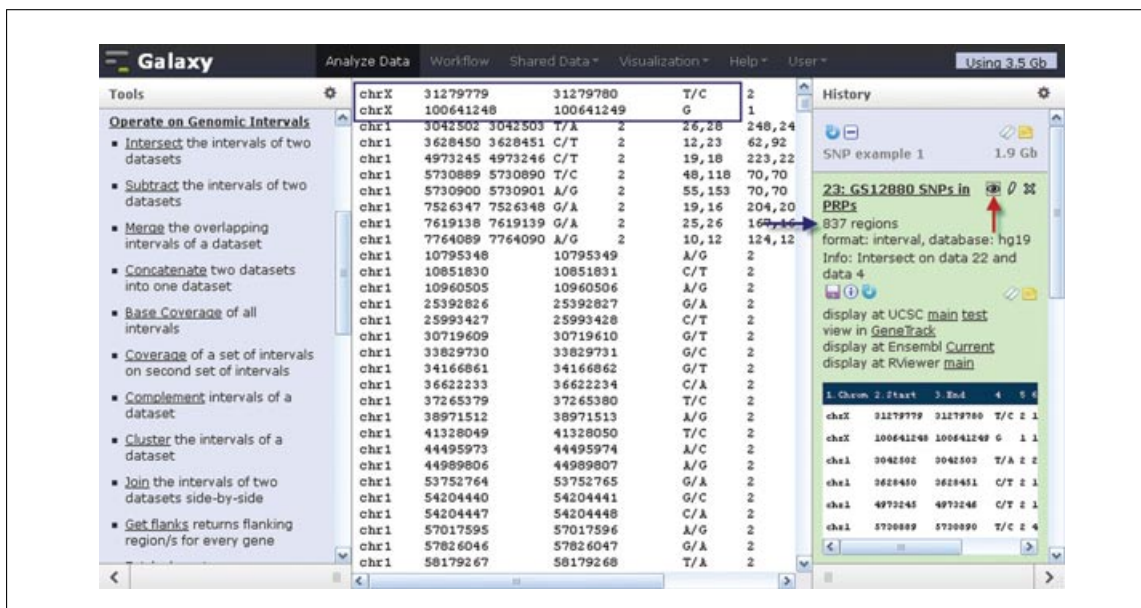
*Intersecting with the PRPs*

2. Exclusion of the SNPs from the 24 healthy individuals is useful for finding both coding and non-coding disease variants, so we will again start with the dataset we prepared in Basic Protocol 1. In the Operate on Genomic Intervals section of the tool panel, select and run the Intersect tool as shown in Figure 15.2.21. An intersection between our input dataset and the PRPs will find the SNPs within the predicted regulatory regions. The output will have the columns from the first dataset, so be sure to specify the SNPs as the first dataset and the PRPs as the second one.
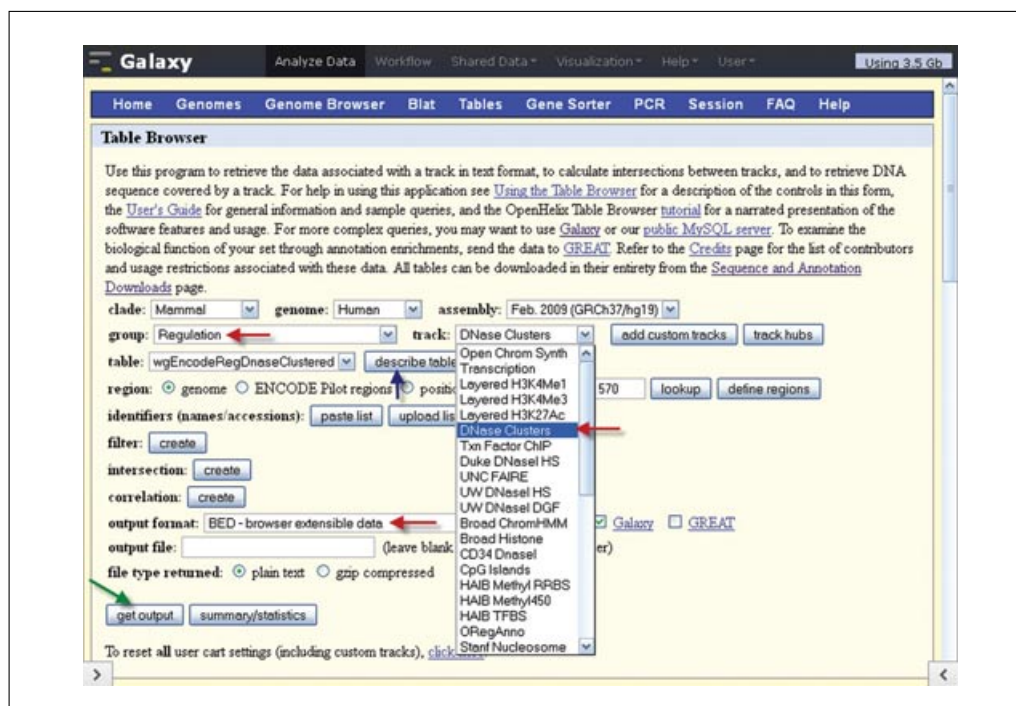
*SNPs in PRPs*

3. There are 837 SNPs in the predicted regulatory regions (Fig. 15.2.22), including two of our non-coding disease SNPs (blue box).



**Figure 15.2.21**    Intersecting with the PRPs. See text for details. For color version of figure go to *http://www.currentprotocols.com/protocol/bi1502*.

**Figure 15.2.22** SNPs in PRPs. See text for details. For color version of figure go to *http://www.currentprotocols.com/protocol/bi1502*.

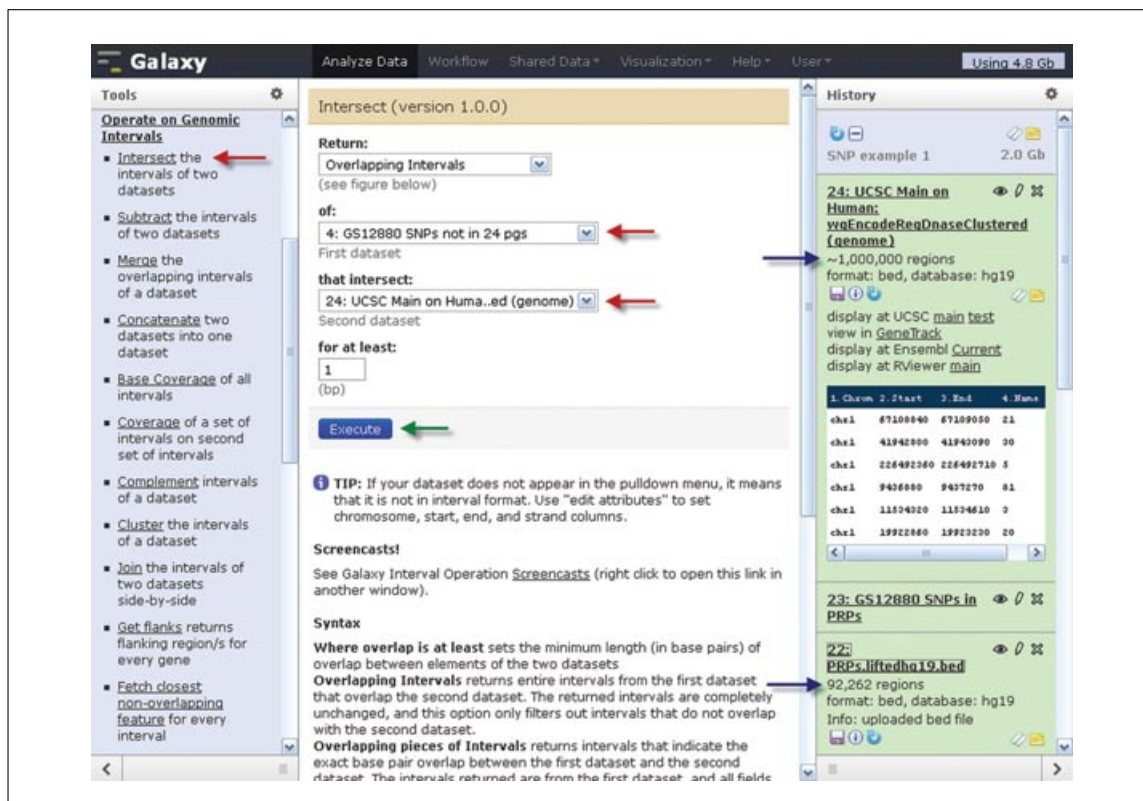

**Figure 15.2.23** DNase hypersensitive sites (HSSs) from ENCODE. See text for details. For color version of figure go to *http://www.currentprotocols.com/protocol/bi1502*.
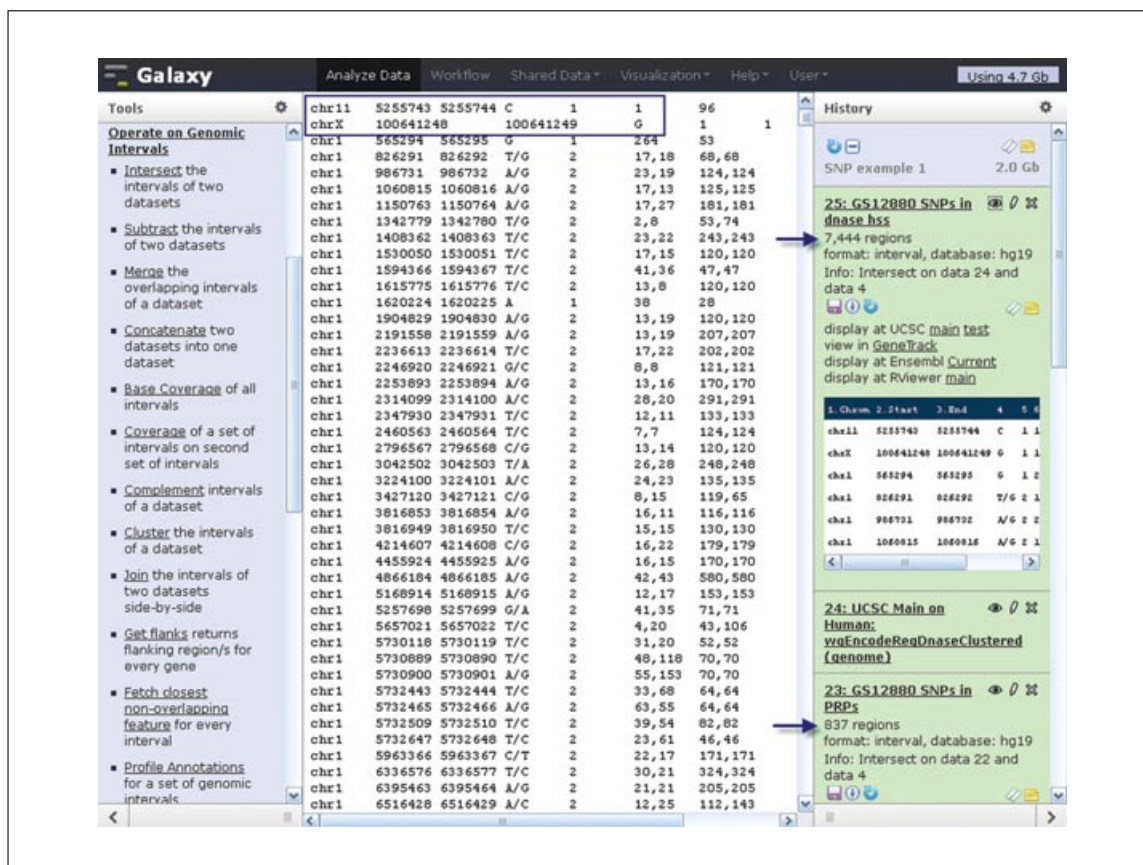
### DNase hypersensitive sites (HSSs) from ENCODE data

4. To obtain data from the ENCODE project, we return to the UCSC Main Table Browser (located in the Get Data section of the tool panel). As before, in order to see the full width of the table browser interface without scrolling, we have temporarily hidden the tool and history panels via the bottom corner arrow buttons (Fig. 15.2.23).

**Figure 15.2.24** Intersecting with the HSSs. See text for details. For color version of figure go to *http://www.currentprotocols.com/protocol/bi1502*.



**Figure 15.2.25** SNPs in HSSs. See text for details. For color version of figure go to *http://www.currentprotocols.com/protocol/bi1502*.

**15.2.21**

**Figure 15.2.26** PhyloP. See text for details. For color version of figure go to *http://www.currentprotocols.com/protocol/bi1502*.



**Figure 15.2.27** Distribution of phyloP scores. See text for details. For color version of figure go to *http://www.currentprotocols.com/protocol/bi1502*.

5. Within the group Regulation there is a large list of tracks (drop-down box), many from ENCODE. For this example we have chosen DNase Clusters, which has only one table available: "wgEncodeRegDnaseClustered." The default format for sending data to Galaxy is BED; this is typically the best format for tables with genomic positions, since it allows doing many genomic operations in Galaxy. Once the options are selected, click the "get output" button. For BED output, an intermediate page then appears with options for how to return the intervals. Just leave these at their defaults and click the "Send query to Galaxy" button.

*In general, when choosing a table, clicking on the "describe table schema" button (blue arrow) will tell you more about the table you are selecting, as well as provide a description of what data and methods were used for the track as a whole.*

### Intersecting with the HSSs

6. The searches we are conducting now are more or less independent, looking for SNPs that affect a variety of region types and functional mechanisms, so rather than

**Figure 15.2.28**   Histogram. See text for details. For color version of figure go to *http://www.currentprotocols. com/protocol/bi1502.*

continuing from the PRP results, we will go back and start yet again with our input dataset from Part A. Once again we are using the Intersect tool, this time to find the SNPs within the DNase hypersensitive sites (Fig. 15.2.24). Note that there are a lot more regions in this set than there were in the PRPs (blue arrows).

### SNPs in HSSs

7. As expected, there are many more results from this intersection (blue arrows in Fig. 15.2.25). Again, two of our disease SNPs are found. One was also found via the PRPs; the other is a non-coding SNP missed by that method.
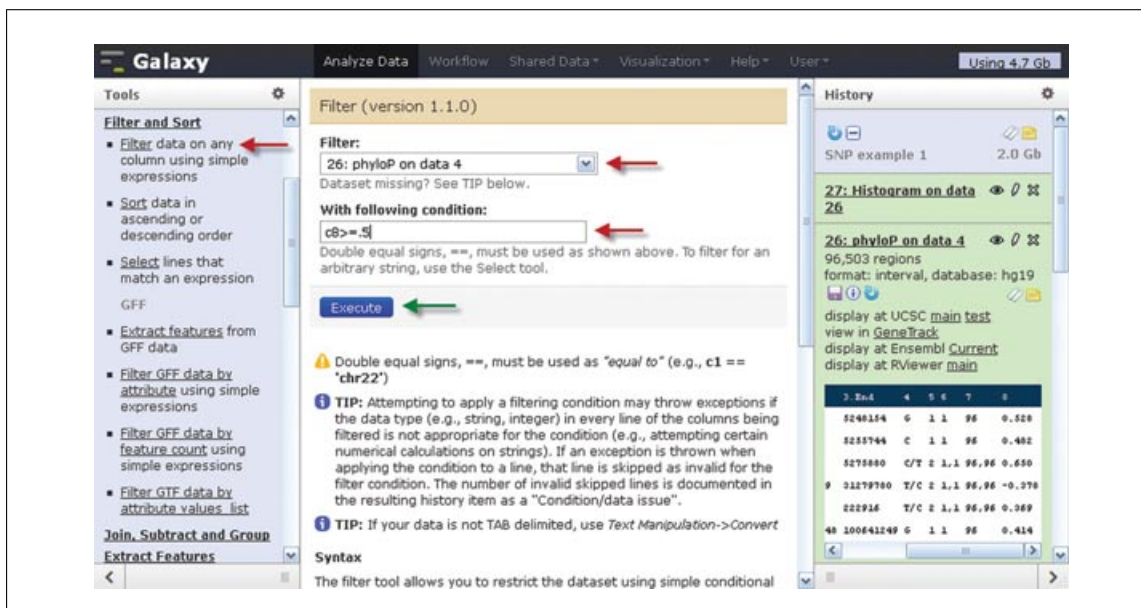
### PhyloP

8. Next, we will look at conservation between species using the PhyloP tool, which looks up a precomputed phyloP score (Siepel et al., 2006) for each SNP position. This tool is located in the Phenotype Association section (Fig. 15.2.26). Select the same input dataset from Basic Protocol 1 that we have been using, and then click Execute.

### Distribution of phyloP scores

9. To use the phyloP scores, it is helpful to know the distribution of those scores. For this we go to the Graph/Display Data section and select the Histogram tool (Fig. 15.2.27). In the history panel, we can see what column the scores are in, and a few of the values (blue box). The PhyloP tool appends a column to its input dataset, and as expected the last column is new and contains numbers. Select the dataset with the phyloP scores and indicate the column that has the score. Entering a large number of bars gives a finer resolution in the output. Adding a descriptive title and label are not required, but are helpful when coming back to this history later. Click Execute to generate the graph.

**Understanding Human Variation**

**15.2.23**

**Figure 15.2.29** Filtering the SNPs based on phyloP score. See text for details. For color version of figure go to *http://www.currentprotocols.com/protocol/bi1502*.



**Figure 15.2.30** Highly conserved SNPs. See text for details. For color version of figure go to *http://www. currentprotocols.com/protocol/bi1502*.

### Histogram

10. Clicking on the eye icon in the history displays the histogram in the center panel (Fig. 15.2.28). Here we have hidden the tool panel to make the graph larger. Any positive phyloP score indicates conservation, but there are many SNPs with scores between 0 and 1. To select the most highly conserved regions, a cutoff of 0.5 seems reasonable.

### Filtering the SNPs based on phyloP score

11. In the Filter and Sort section, choose the Filter tool (Fig. 15.2.29). Select the dataset with the phyloP scores, and enter the condition c8>=.5, which means we want the rows where the score in column 8 is greater than or equal to 0.5.

### *Highly conserved SNPs*

12. This approach has the highest number of results yet (first blue arrow in Fig. 15.2.30), and finds two of our disease SNPs, including one that has not been found up to this point (blue box). The message about skipped lines (second blue arrow) is not a problem; NA in the PhyloP tool's output means that the score was not available, and we want to skip those SNPs anyway.

> *This is the end of the protocol. To narrow the results further, one could intersect the output from the various searches, either pairwise or all together.*

## GUIDELINES FOR UNDERSTANDING RESULTS

It is important to use caution when working with SNP data, as there are always some errors in the sequencing and SNP-calling processes. Even an error rate as low as 0.001% (as reported by Complete Genomics) can produce twenty to thirty errors among the two to three million SNPs typically found in an individual, and in practice error rates often seem much higher (e.g., $\sim$1%), especially for some technologies.

Furthermore, genotype-to-phenotype associations are rarely straightforward. Phenotype association does not necessarily mean phenotype causation; other factors also influence the phenotype, such as other genes, lifestyle, and the environment. For complex diseases, associations determined using genome-wide association studies often reflect just an increase or decrease in risk. The SNP listed may not even be the true phenotype-affecting SNP, but just one that happens to be in linkage disequilibrium with the true SNP. Also, keep in mind that the predictions made for coding SNPs by programs such as PolyPhen-2 and SIFT are only predicting the effect of the SNP on the protein, not on the individual as a whole.

Conversely, many of these caveats apply not only to the SNPs found, but also to the so-called "benign" ones we eliminated. Just because an allele is reported in an apparently healthy individual does not guarantee that it never contributes to problems in other people, or in the same person later on, or even that the person tested actually has that allele.

In summary, results from these protocols should be considered only as rough hints or clues to suggest candidate SNPs potentially worthy of further investigation.

## COMMENTARY

### Background Information

Data on human genetic variation (e.g., from genotyping, whole-genome sequencing, genome-wide association studies, and mapping of epigenetic features) have the potential to advance medical applications at an unprecedented rate, but only if bioinformatics solutions are provided to bench scientists and clinicians to enable them to find biological insights in the flood of data.

Currently, most software tools, databases, and other resources for these analyses are either not publicly available, or are scattered among many sites and require considerable effort and expertise to install and operate. In several cases, projects to produce open-source resources have published papers, but the resulting tools have not yet been made available.

Galaxy serves as an effective platform for assembling a wide variety of tools into an integrated, easily accessible system that offers one-stop convenience, no need for users to install specialized software, and most importantly, a consistent interface that does not require programming expertise to use. The Phenotype Association section provides access to many tools that are useful for analyzing biomedical applications of genetic variation data, only a few of which are discussed here.

### Critical Parameters and Troubleshooting

Be sure to read the description below each tool's input form, especially with regard to the input and output formats.

For each dataset, Galaxy keeps track of the data format and the genome assembly (for datasets containing genomic positions). If a tool takes multiple datasets as input, Galaxy will generally require that they be based on the same assembly (except for tools that are specifically designed to work across species or builds). It will also only allow you to select input datasets that are in a format compatible with the tool. Therefore, if the dataset you want to use as input for a tool is in your history but is not appearing in the list to select from, check the assembly (also known as "database") and format of the dataset to make sure they are correct. There are tools for converting between formats when needed.

Note that if Galaxy ever guesses the type of a dataset incorrectly, you can fix that by clicking on its pencil icon in the history and using the "Change data type" form at the bottom of the page that appears. This is different from actually converting the data from one format to another.

For further assistance, please consult the Critical Parameters and Troubleshooting section of *UNIT 10.5*.

## Literature Cited

Adzhubei, I.A., Schmidt, S., Peshkin, L., Ramensky, V.E., Gerasimova, A., Bork, P., Kondrashov, A.S., and Sunyaev, S.R. 2010. A method and server for predicting damaging missense mutations. *Nat. Methods* 7:248-249.

Blankenberg, D., Von Kuster, G., Coraor, N., Ananda, G., Lazarus, R., Mangan, M., Nekrutenko, A., and Taylor, J. 2010. Galaxy: A web-based genome analysis tool for experimentalists. *Curr. Protoc. Mol. Biol.* 89:19.10.1-19.10.21.

Davis, A.P., Murphy, C.G., Saraceni-Richards, C.A., Rosenstein, M.C., Wiegers, T.C., and Mattingly, C.J. 2009. Comparative Toxicogenomics Database: A knowledgebase and discovery tool for chemical.gene.disease networks. *Nucleic Acids Res.* 37:D786-D792.

Drmanac, R., Sparks, A.B., Callow, J.M., Halpern, A.L., Burns, N.L., Kermani, B.G., Carnevali, P., Nazarenko, I., Nilsen, G.B., Yeung, G, Dahl, F., Fernandez, A., Staker, B., Pant, K.P., Baccash, J., Borcherding, A.P., Brownley, A., Cedeno, R., Chen, L., Chernikoff, D., Cheung, A., Chirita, R., Curson, B., Ebert, J.C., Hacker, C.R., Hartlage, R., Hauser, B., Huang, S., Jiang, Y., Karpinchyk, V., Koenig, M., Kong, C., Landers, T., Le, C., Liu, J., McBride, C.E., Morenzoni, M., Morey, R.E., Mutch, K., Perazich, H., Perry, K., Peters, B.A., Peterson, J., Pethiyagoda, C.L., Pothuraju, K., Richter, C., Rosenbaum, A.M., Roy, S., Shafto, J., Sharanhovich, U., Shannon, K.W., Sheppy, C.G., Sun, M., Thakuria, J.V., Tran, A., Vu, D., Zaranek, A.W., Wu, X., Drmanac, S., Oliphant, A.R., Banyai, W.C., Martin, B., Ballinger, D.G., Church, G.M., and Reid, C.A. 2009. Human genome sequencing using unchained base reads on self-assembling DNA nanoarrays. *Science* 327:78-81.

Ferretti, V., Poitras, C., Bergeron, D., Coulombe, B., Robert, F., and Blanchette, M. 2007. PReMod: A database of genome-wide mammalian *cis*-regulatory module predictions. *Nucleic Acids Res.* 35:D122-D126.

Giardine, B., Riemer, C., Hardison, R.C., Burhans, R., Elnitski, L., Shah, P., Zhang, Y., Blankenberg, D., Albert, I., Taylor, J., Miller, W., Kent, W.J., and Nekrutenko, A. 2005. Galaxy: A platform for interactive large-scale genome analysis. *Genome Res.* 15:1451-1455.

Giardine, B., Riemer, C., Hefferon, T., Thomas, D., Hsu, F., Zielenski, J., Sang, Y., Elnitski, L., Cutting, G., Trumbower, H., Kern, A., Kuhn, R., Patrinos, G.P., Hughes, J., Higgs, D., Chui, D., Scriver, C., Phommarinh, M., Patnaik, S.K., Blumenfeld, O., Gottlieb, B., Vihinen, M., Väliaho, J., Kent, J., Miller, W., and Hardison, R.C. 2007. PhenCode: Connecting ENCODE data with mutations and phenotype. *Hum. Mutat.* 28:554-562.

Goecks, J., Nekrutenko, A., Taylor, J.; Galaxy Team. 2010. Galaxy: A comprehensive approach for supporting accessible, reproducible, and transparent computational research in the life sciences. *Genome Biol.* 11:R86.

Huang, D.W., Sherman, B.T., and Lempicki, R.A. 2009. Systematic and integrative analysis of large gene lists using DAVID bioinformatics resources. *Nat. Protoc.* 4:44-57.

Karolchik, D., Hinrichs, A.S., Furey, T.S., Roskin, K.M., Sugnet, C.W., Haussler, D., and Kent, W.J. 2004. The UCSC Table Browser data retrieval tool. *Nucleic Acids Res.* 32:D493-D496.

Kumar, P., Henikoff, S., and Ng, P.C. 2009. Predicting the effects of coding non-synonymous variants on protein function using the SIFT algorithm. *Nat. Protoc.* 4:1073-1081.

Reimand, J., Kull, M., Peterson, H., Hansen, J., and Vilo, J. 2007. g:Profiler: A web-based toolset for functional profiling of gene lists from large-scale experiments. *Nucleic Acids Res.* 35:W193-W200.

Seal, R.L., Gordon, S.M., Lush, M.J., Wright, M.W., and Bruford, E.A. 2011. genenames.org: The HGNC resources in 2011. *Nucleic Acids Res.* 39:D514-519.

Siepel, A., Pollard, K.S., and Haussler, D. 2006. New methods for detecting lineage-specific selection. *In* Proceedings of the 10th International Conference on Research in Computational Molecular Biology (RECOMB 2006), pp. 190-205, Venice, Italy.

Taylor, J., Tyekucheva, S., King, D.C., Hardison, R.C., Miller, W., and Chiaromonte, F. 2006. ESPERR: Learning strong and weak signals in genomic sequence alignments to identify functional elements. *Genome Res.* 16:1596-1604.

## Internet Resources

http://galaxyproject.org
*The main public instance of Galaxy.*

http://phencode.bx.psu.edu
*A collection of human phenotype-associated SNPs from Locus-Specific Databases.*

http://www.bx.psu.edu/miller_lab/docs/galaxy_phen_assoc/tutorial/
*A version of this tutorial in HTML format.*

http://genome.ucsc.edu/FAQ/FAQformat.html
*Descriptions of file formats used by the UCSC Table Browser.*

## Supplementary File

http://www.currentprotocols.com/protocol/bi1502
*This is an alternate URL to access the file "test.masterVar.gz" cited in Basic Protocol 1, Necessary Resources, Files on page 15.2.3.*

**Understanding Human Variation**

**15.2.27**