# Tutorial

## Batching of Multi-Input Workflows

November 21, 2017

Sample to Insight

# Batching of Multi-Input Workflows

This tutorial will demonstrate how workflows that have more than one variable input can be run in batch mode. This is the case for many Ready-to-Use workflows available in the toolbox of Biomedical Genomics Workbench, including all workflows that analyze a trio or a family of fours and workflows that compare tumor to normal samples. In this tutorial we will use the "Identify Causal Inherited Variants In Trio" workflow to illustrate how multiple families' genomes can be analyzed in batches. Note that you must be working with the Biomedical Genomics Workbench 3.0 or higher to be able to batch workflows with multiples inputs.

## Download and import data

To allow you to run this tutorial within minutes, we have generated for each family member a subset of whole exome data containing only chromosome 21. We also introduced a clinically relevant variant in the proband and one of its parent. The Biomedical Genomics Workbench offers a reference data set specific to chromosome 21 and available from the Data Manager.

The data set you will download includes:

- The reads from the father, mother and proband for each of the two families A and B.

- A target regions file called *SeqCap_EZ_Exome_v2_BED_CHR21*. Please note that when working with your own data, you can obtain the relevant target region tracks from the vendor of the amplicon or hybridization kit.

- An Excel spreadsheet with information about the samples, such as which reads belong to which family, and which family members are affected by a particular disease.

Go through the following steps to download and import the data into the Workbench.

1. Download the sample data from our website: `http://resources.qiagenbioinformatics.com//testdata/batching_workflows_tutorial.zip`.

2. Unzip the zip file on your local machine.

3. Start the *Biomedical Genomics Workbench*.

4. Launch **File** | **Import** (📥) | **Standard Import** (📥) and choose the subfolder called "reads and target region" that contains the sequencing reads and the target sequence via the toolbar. Leave the import type set to Automatic and save the imported data in a folder you create in the Navigation Area called "tutorial".

5. You do not need to import the Excel file at this point.

Once the data has been downloaded and imported, you should see the folder and files in the Navigation Area as shown in figure 1.

## Data management configuration

As the data used in this tutorial is from chromosome 21 only, you will need to specify a Reference Data Set that contains only chromosome 21. This Reference Data Set can be found and selected in the Data Manager.
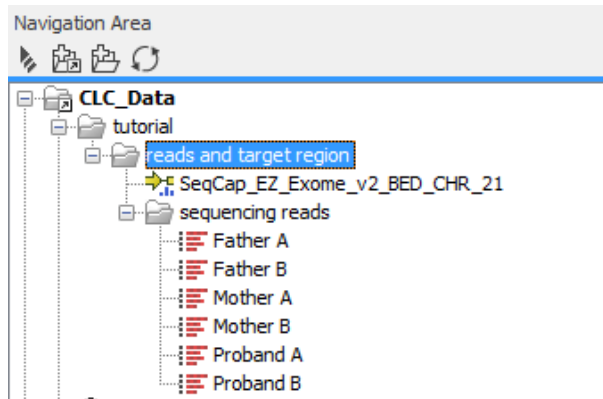
Figure 1: *The data set has been imported.*

1. Click on the button labeled **Data Management** (📂) in the top right corner of the Workbench.

2. In the top right corner of the Data Manager, you can see whether you are managing the reference data locally or on a server (if you have one). Make sure you are downloading the Reference Data Set in the location where you will run the workflow, i.e., "locally" if you want to run the workflow on your machine, and "on the server" if you intend to run the analyses on a server or on a grid.

3. Click on Tutorial Reference Data Sets and choose "Batching of Multi-Input Workflows".

4. In the new window (figure 2) click on the **Download** button to import in your CLC_References folder the elements that were not yet downloaded (marked with an (⊕) icon). After downloading, check that all elements are marked with an (✓) icon.
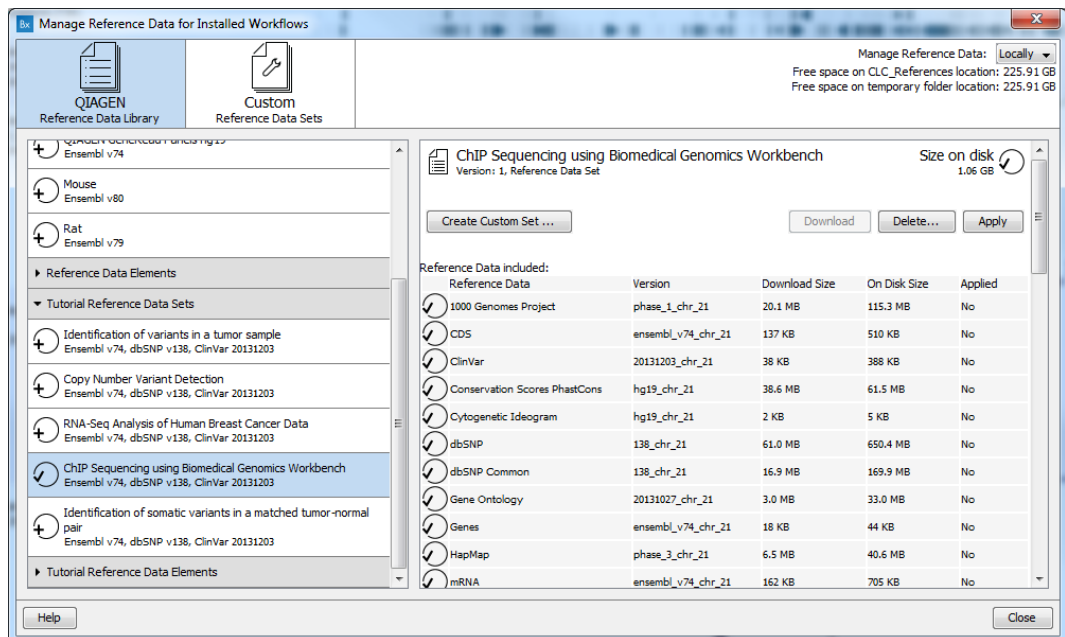


Figure 2: *Apply the Reference Data Set called "Batching of Multi-Input Workflows" in the Data managemet.*

5. Once all elements have been downloaded, click the **Apply** button in the top right corner, and then click **Close** to exit the Data Manager.

Note: remember to change back the set of references applied once you are done with this tutorial by following the instructions given at the end of this tutorial.

**Running a Ready-to-Use workflow in batch mode**

In this tutorial we will analyze two families at once using the batch option using a workflow containing a Trio. The reads for each family member are now available in the Navigation Area, but to be able to run both families at the same time, you need to provide additional information (organized in an Excel spreadsheet) about the origin and the grouping of the reads: which family member do the samples belong to, who is affected by the disease, etc.

To find the workflow, go to:

> **Toolbox | Ready-to-Use Workflows | Whole Exome Sequencing (📲) | Hereditary Disease (📂) | Identify Casual Inherited Variants In Trio (WES)**

1. Right click on the **Identify Casual Inherited Variants In Trio (WES)** tool. This will open a menu as shown on figure 3. Choose "Run in Batch Mode".
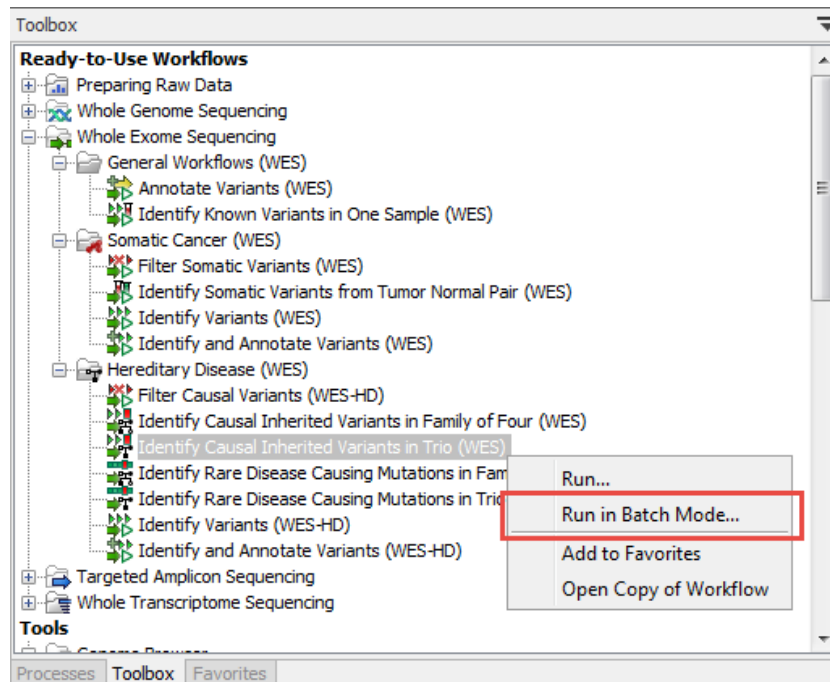


Figure 3: *Right clickling the Identify Casual Inherited Variants In Trio (WES) tool will open a window giving the posibility to choose Run in Batch Mode.*

2. Depending on your local setup you might be asked to select where to run the workflow. Choose the appropriate option for you and click **Next**.

3. In the window called "Import metadata":

- In the "Spreadsheet with sample information" section, select the Excel spreadsheet with the sample information that you downloaded on your computer earlier.

- In the "Data to analyze" section select the folder containing the reads for the two families as shown on figure **??**. Note that to analyze multiple families in one batch run, the reads must be stored together in one or several folders, as it is only possible to select folders, and not individual files.

- In the "Data association" section, all samples should now be marked with a green check mark and the wizard window should look like figure 4. Click **Next**.
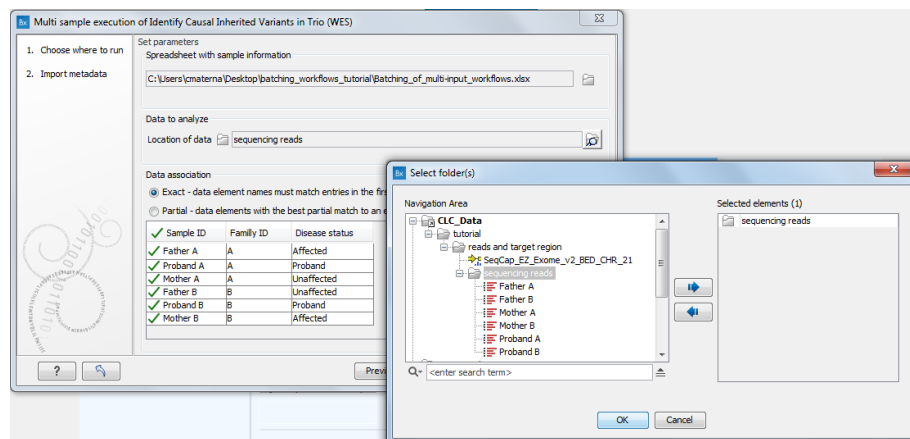


Figure 4: *The Import metadata wizard window after the spreadsheet and the reads were specified and correctly associated.*

4. In the next wizard window (figure 5):

- In the "Select grouping parameters and analysis inputs" section, choose for "Group by" the option **Family ID** from the dropdown menu, as we are here grouping the samples by families (one batch is one family). For the "Type", choose the option **Disease status**, as the "Type" indicates who is the affected parent, the unaffected parent and the proband.

- In the "Sample columns" section, select successively "Proband", "Affected" and "Unaffected" in the relevant fields.

- Finally select *SeqCap_EZ_Exome_v2_BED_CHR21* as target region. The wizard window should now look like figure 5.

Click **Next**.

5. In the "Fixed Ploidy Variant Detection" windows, keep "Required variant probability" and "Minimum frequency" at their default values of 50.0 and 10.0 , but set **Minimum coverage** and **Minimum count** to 40 for all family members independently, i.e, in the three successive windows specific to proband, affected parent and unaffected parent (figure 6). We usually recommend to use default settings but have chosen here more stringent parameters for demonstration purposes.

6. In the "QC for Target Sequence" wizard window leave the parameters as default with Minimum coverage set to 30 and leave the options "Ignore non-specific matches" and "Ignore broken pairs" unchecked (figure 7).
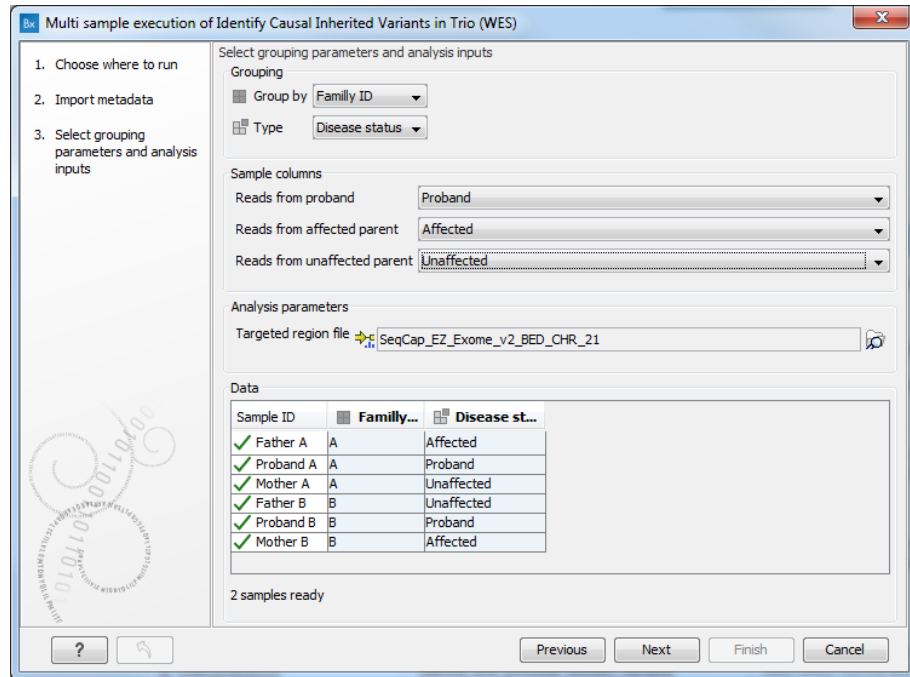
Figure 5: *The Multi sample execution of Identify Casual Inherited Variants in Trio (WES) after the parameters have been set.*
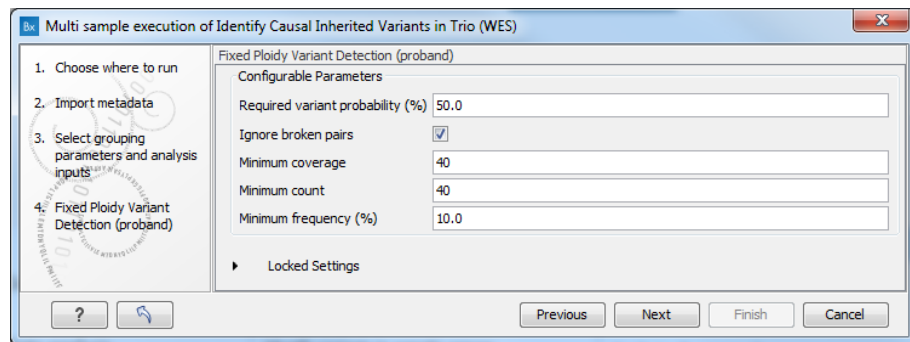


Figure 6: *Set Minimum count and Minimum coverage to 40 for the "Fixed Ploidy Variant Detection" tool.*

7. In the "Remove Variants Found in HapMap" window, leave the 12 elements as they are already selected (figure 8) and click **Next**.

8. Finally, choose to save the results in a separate folder: for example create a new folder called Results.

## Analyzing the results

Once the analysis has completed, the results are organized by batch, with one subfolder per family in the Results folder as can be seen on figure 9. Each folder contains a Genome Browser View that can be opened by double clicking on it from the Navigation Area. The View opens together with a variant table listing the variants identified in the proband. These variants are the putative disease causing variants as they are shared with the affected parent, and not with the
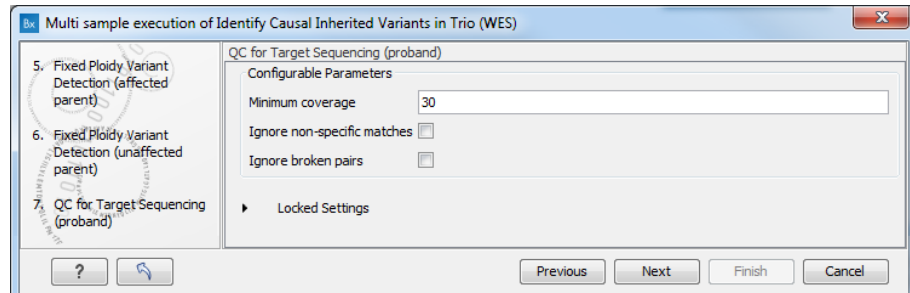
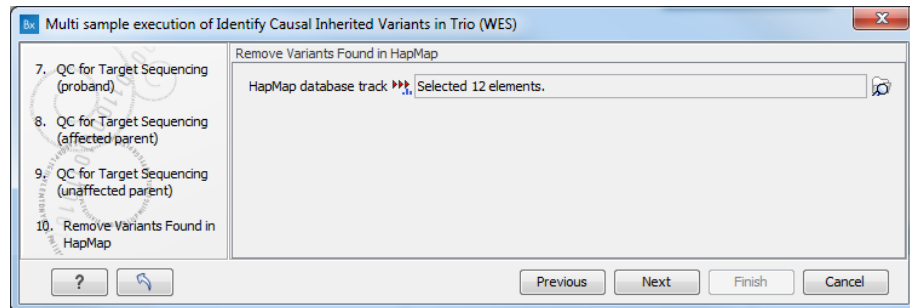Figure 7: *Set the Minimum coverage in the "QC for Target Sequence" wizard window as 30.*



Figure 8: *We willwork with the 12 populations from the hapmap database.*
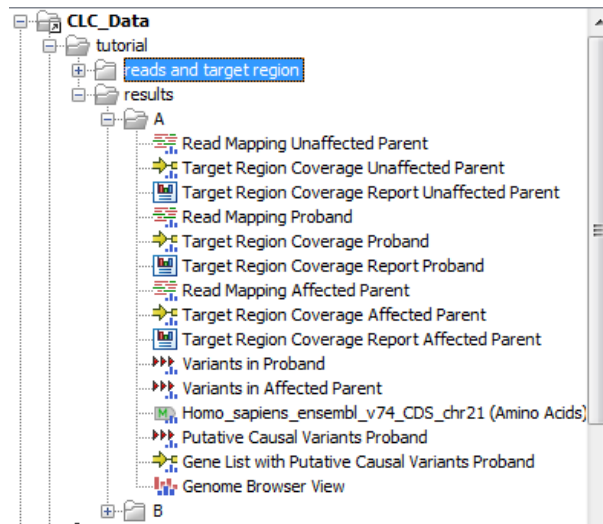
unaffected parent.



Figure 9: *When running a workflow in a batch mode, one folder per batch is generated.*

The Genome Browser View and the table below are linked, so clicking on a variant in the table will display that variant in the Browser View at the nucleotide level, allowing you to look closely at the read mapping and coverage at a particular position. You get the opportunity to compare the reads in the unaffected parent with the ones in the affected parent and the proband (figure 10).

It also allows you to take a look at the amino acids and see if the variant has resulted in any change in the amino acid sequence (figure 11).

The table below also shows information about changes in the allele compared to the reference,
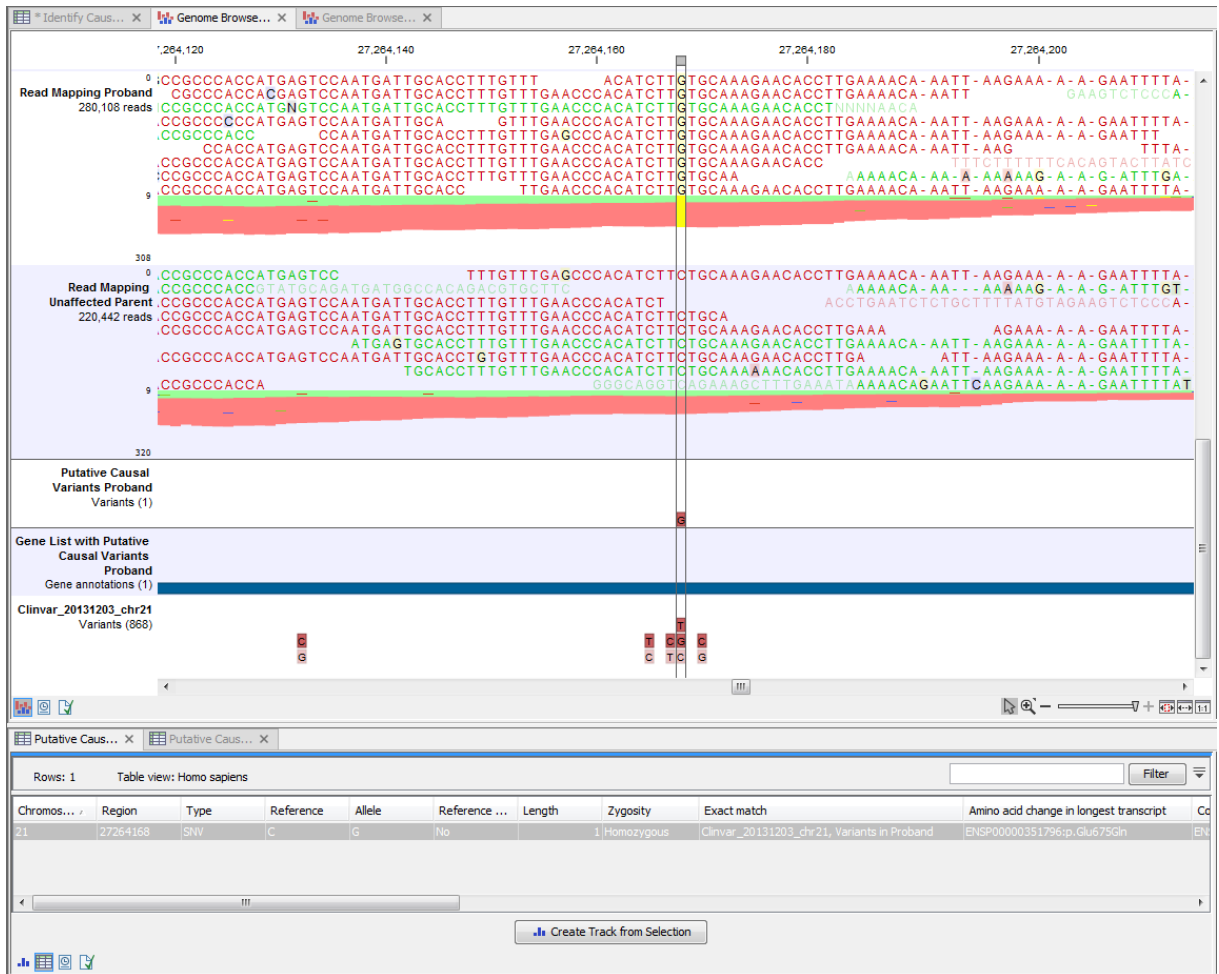
Figure 10: *When viewing the variant at the nucleotide level, it is possible to compare the reads in the different family members and get an overview of whether the variant is already known in a particular database.*

amino acid changes, information about which gene carries the change, and suggest which disease the change might be involved in. In this tutorial we have identified one mutation per family: for family A, the mutation is in the amyloid precursor protein gene and causing hereditary cerebral hemorrhage with amyloidosis - Dutch type; for family B, the variant is in the alphaA-crystallin gene CRYAA and linked to inherited cataract-microcornea.

## Data Management reminder

Before you leave this tutorial, remember to change the applied Reference Data Set back to a full genome version.

1. Open the **Data Management** tool by clicking on the button with that name in the top right corner of the Workbench.

2. If you have a server, choose whether you want to apply the Reference dataset locally or on the server.
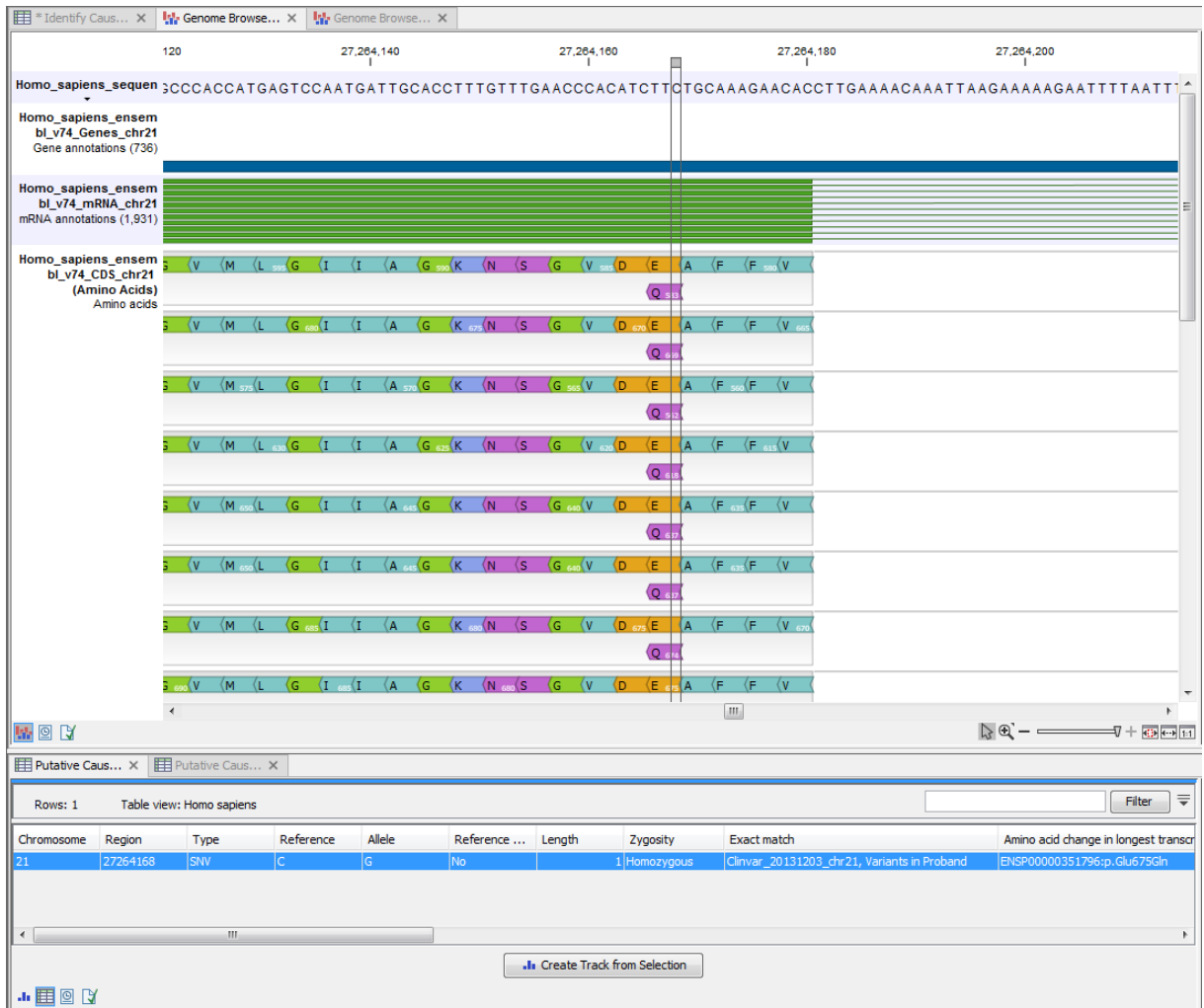
Figure 11: *The amino acids trackin the Genome Browser View shows potential amino acids changes that have occurred as the result of the variant.*

3.  Select the **hg19** or **hg38** Reference Data Set.

4.  Click on the button labeled **Apply** before clicking on the button labeled **Close**.