



Tutorial

Comparative Analysis of Three Bovine Genomes

September 10, 2019

— Sample to Insight —

Comparative Analysis of Three Bovine Genomes

This tutorial takes you through some of the tools for identifying species-specific variants. As an example we use *Bos taurus* (cow), *Bos indicus* (zebu), and *Bison bison* (bison) data sets consisting of single GA and GAI exome sequencing short reads from an Illumina sequencer. For this tutorial, you will need to use CLC Genomics Workbench Version 9.5 or higher.

The goal is to identify zebu-specific (*Bos indicus*) exome variations that are not present in the bison (*Bison bison*) or bovine genome (*Bos taurus*) [Cosart et al., 2011]. Furthermore, we attempt to link the variations to altered pathways using the gene ontology tool that is built into the CLC Genomics Workbench.

Downloading the reference data

Download the bovine genome with annotations using the Reference Data Manager (1) found in the upper right corner of the Workbench (figure 1). Under the **Download Genome** tab (2), select the **Bos taurus - Hereford** data (3). Choose to "Download genome sequence" (4) and check the "Genome Annotations" item (5) to get corresponding annotation tracks. Click on **Download** (6). You can check the Download process at the bottom of the wizard.

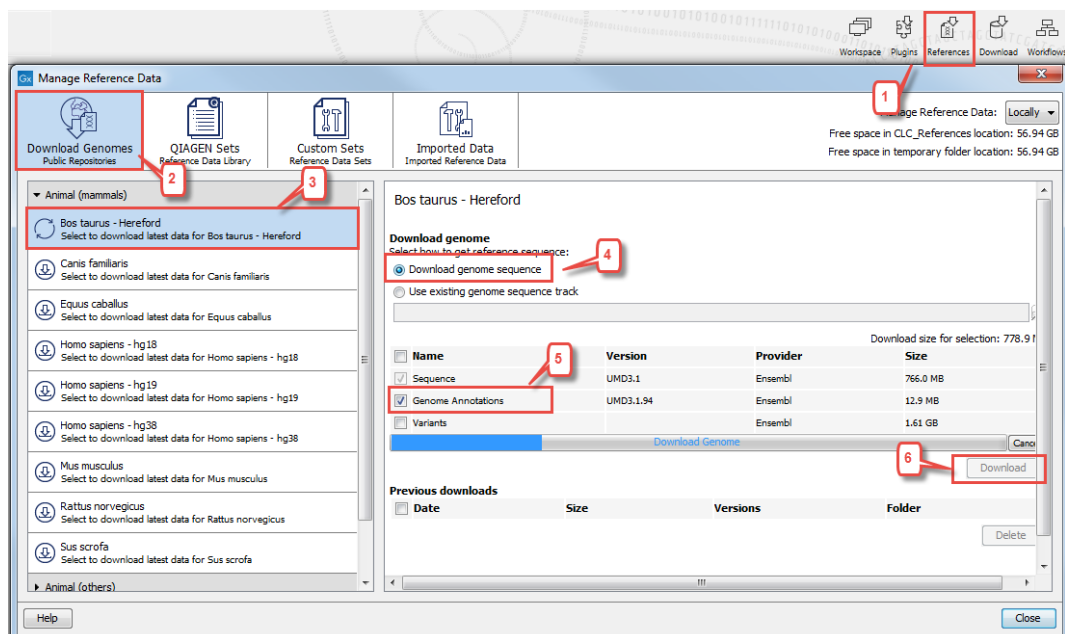


Figure 1: The Reference Data Manager is found in the upper right corner of the Workbench.

The reference data sequence and annotations tracks are now saved in the CLC_References | Genomes | Bos_taurus_Hereford folder accessible from the Navigation Area.

Importing reads

1. Download the read files using the tool:

Download | Search for Reads in SRA 

2. This will bring up the dialog shown in figure 2. Search for one of the runs with the following accession number SRR307349.

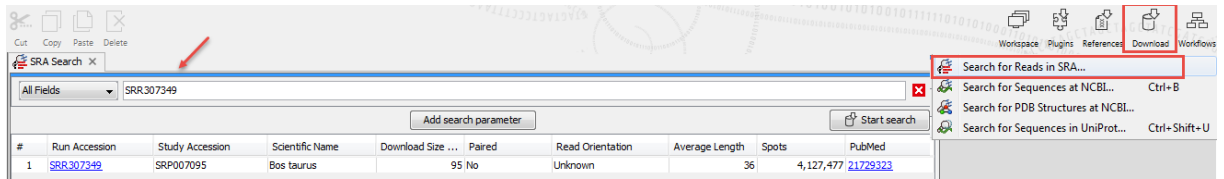


Figure 2: Search for a particular run using the Search for Reads in SRA tool.

3. The search output one run. To find the other reads generated during this study, right click on the row of the search output table and choose the option "search for..." and select SRP007095, i.e., the accession number of the study (figure 3).

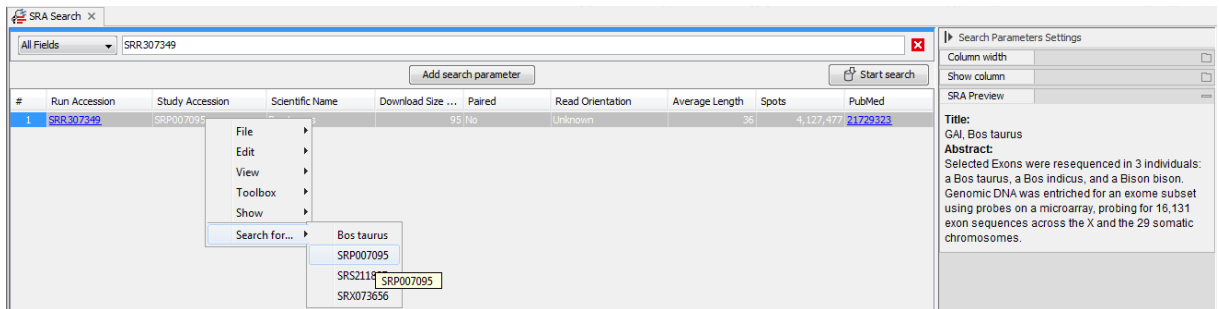


Figure 3: Search for more runs from the same experiment.

4. You have now 6 different runs in the search output table. To minimize the tutorial download and analysis time, select the 3 runs that have the lowest Download Size and click on Download Reads and Metadata (figure 4).

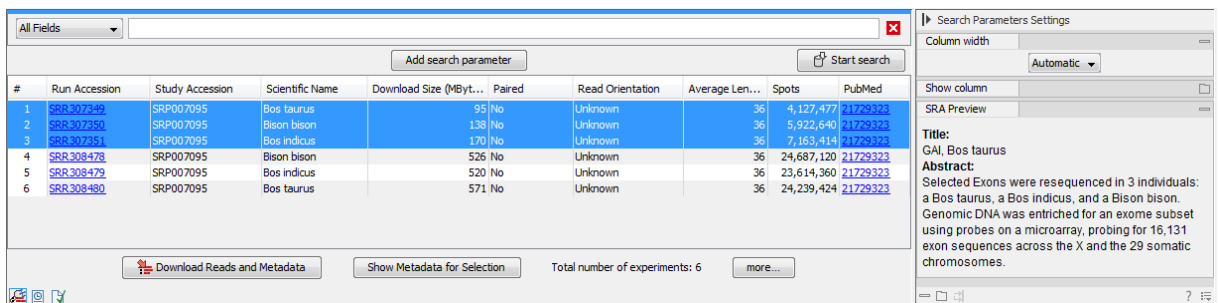


Figure 4: Download in one click the runs and associated metadata.

5. In the download dialog that pops up, choose to discard read names and quality scores.
6. Save the data in the folder you created for this tutorial.
7. For clarity in subsequent analysis, rename the reads with the corresponding species name (right click on the read file in the Navigation Area).

When all data are downloaded you can see the following files in the **Navigation Area** as shown in figure 5.

We are now ready to start analyzing the data.

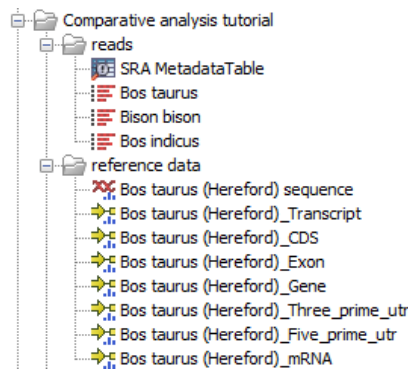


Figure 5: The downloaded reference genome data.

Map the reads and detect variants

The data used in this tutorial do not contain adapters, so we do not need to run a trimming step. Thus, we proceed directly to mapping the reads to the reference.

We will map the reads to the reference genome using the "Batch" mode, which means that the three sets can be run in parallel (simultaneously).

1. To map reads to the reference, go to the toolbox:

Toolbox | Resequencing Analysis | Map Reads to Reference

In the first wizard step (figure 6), check the **Batch** box, then select the folder that holds the reads. Do not forget to check the batch box as it would generate a single read mapping with the reads of all three species. Click **Next**.

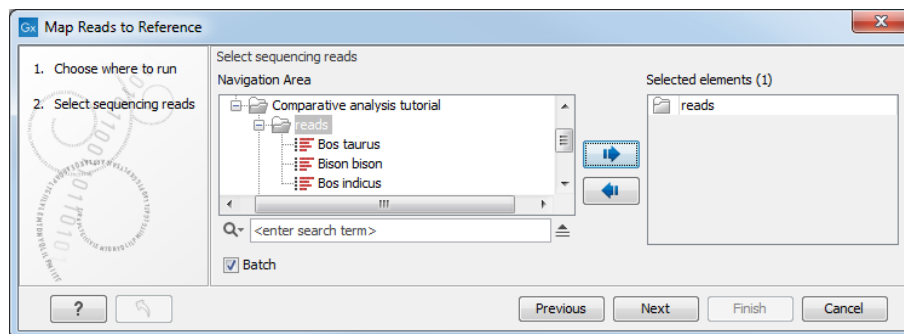



Figure 6: Run the analysis in batch mode to generate one read mapping per species.

2. You should now see the batch overview window (figure 7). In this wizard step you can see the reads that were contained in the selected folder. In this case we will use all the reads in the analysis so you can just click **Next**.
3. In the next wizard window (figure 8), click on the browse icon  under **References** to select the *Bos taurus* (Hereford) reference sequence. Press OK. Leave the Reference masking parameters as set by default to "No masking" and click **Next**.
4. In the next wizard step, you can specify the **Mapping options**. We will use the default settings and click **Next**.
5. Select the options "Create reads track", "Create report" and "Save in specified location" using the option to "Create subfolders per batch unit". Specify where you want to save the data, for example a new folder called "Analysis" and click on the button labeled **Finish**.

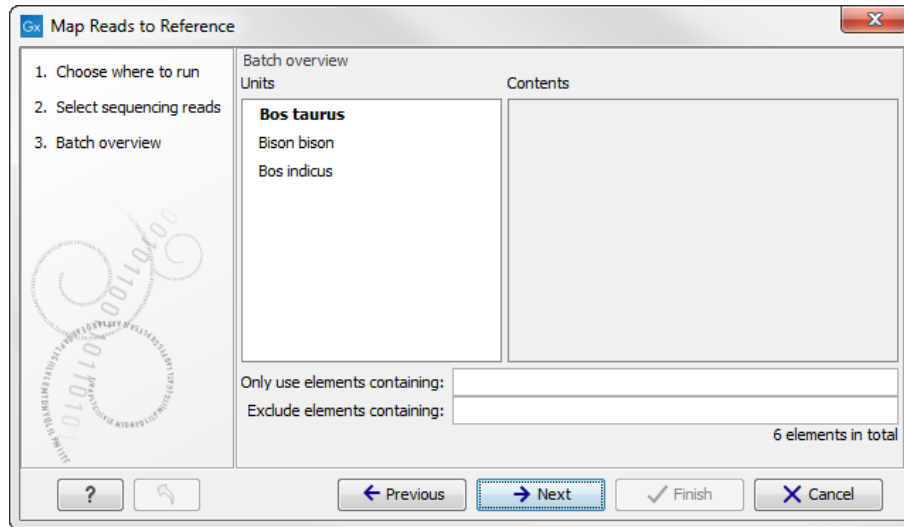


Figure 7: Failing to see this dialog means you forgot to check the batch option previously!

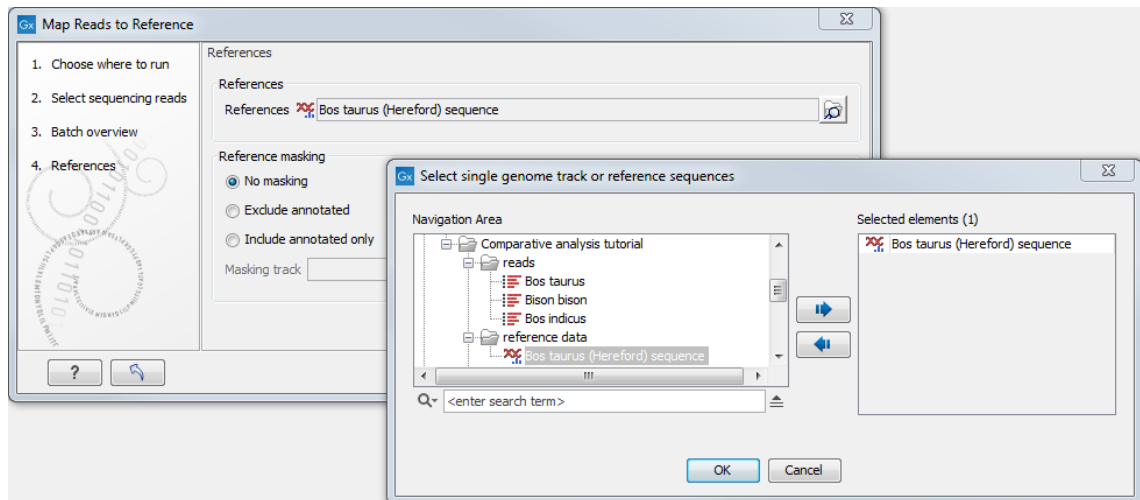


Figure 8: Select the *Bos taurus (Hereford)* reference sequence.

Open the mapping summary reports from the **Navigation Area**. 100% of reads are mapped successfully in all three analyses. We can then proceed to a local realignment of the reads to improve the quality of the read mappings we have just generated.

1. To run the local realignment tool go to the toolbox:

Toolbox | Resequencing Analysis (📁) | Local Realignment

2. Again we choose to run the analysis in **batch mode**, which allows us to analyze all samples in parallel. Select all the read mapping tracks by selecting the top folder as shown in figure 9 and click **Next**.
3. The next dialog in which you can see the content of the selected folder confirms that you are working in batch mode. Click **Next**.
4. In the next wizard step, use the default settings (we will run the local realignment without a guidance track). Click on the button labeled **Next**.

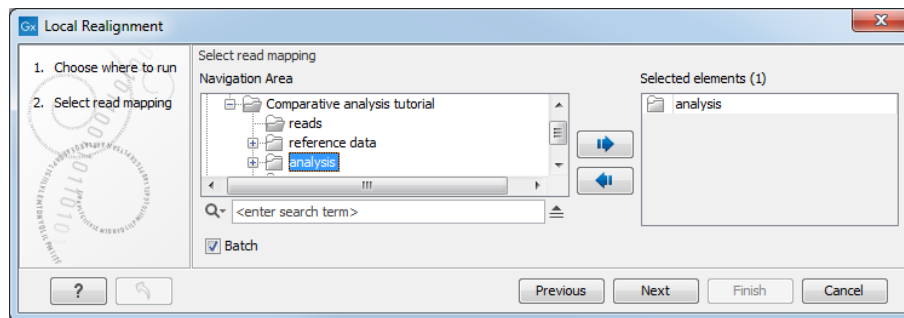


Figure 9: Run local realignment for all species at the same time by using the batch function.

5. Click on **Create reads track** (see figure 10) and choose to **Save in input folder**. Click on **Finish**.

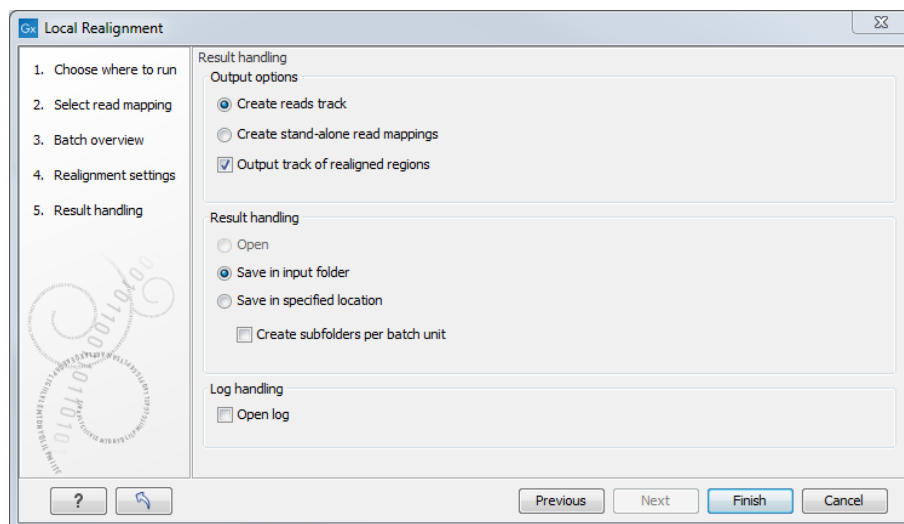


Figure 10: Use the default settings and choose to save the results.

We are now ready to identify the variants. Note that we are generating here a variant track for all three species, when in fact we will use only the *Bos indicus* variant track in subsequent analyses.

1. Go to the toolbox:

Toolbox | Resequencing Analysis  | **Variant Detection**  | **Basic Variant Detection** 

2. Using the batch function, select the top folder as in figure 11.
3. Click through the next wizard steps (by clicking on the button labeled Next) using the default settings. If you have changed some of the settings, you can always go back to the default settings by pressing the **Reset** button in the lower left corner, which will return parameters to default settings.
4. In the last wizard step choose to **Create reads track** and **Save in input folder** before clicking on **Finish** (figure 12).

Take a look at the results, which should be organized as shown in figure 13.

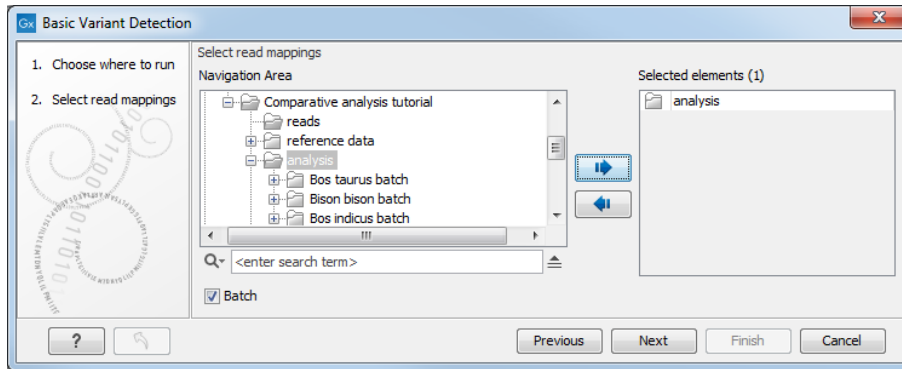


Figure 11: Select the folder that holds the locally realigned reads.

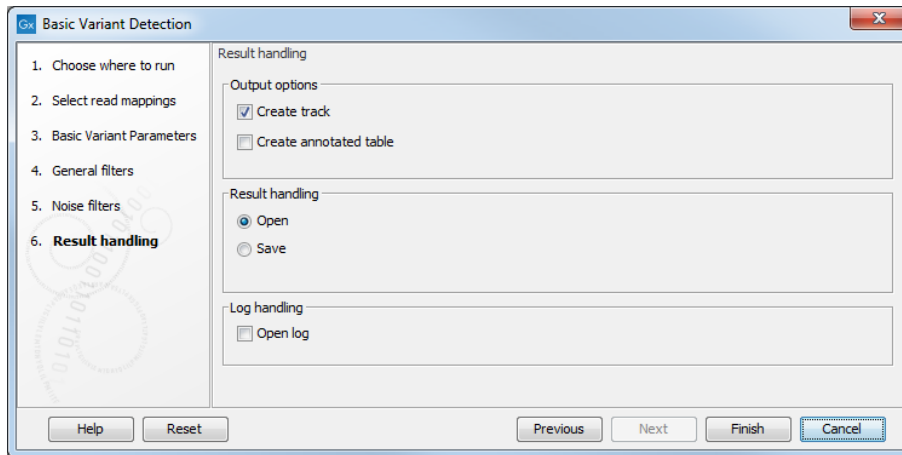


Figure 12: Create a reads track and save the results in their respective input folders.

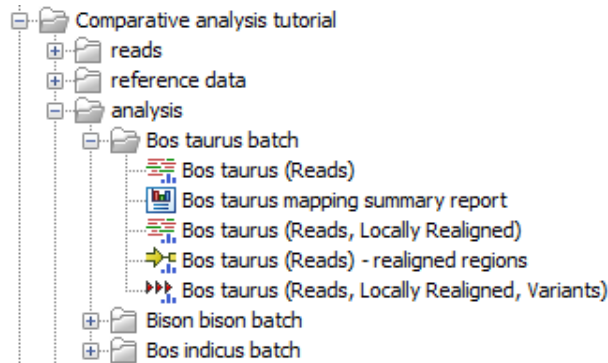


Figure 13: Example of how all outputs have been saved to the batch specific folder created during the first batch analysis.

Let's take a look at the one of the variant tracks. Open the **Bos taurus (Reads) - locally realigned (Variants)** track by double-clicking on the name in the **Navigation Area**. Use the lower left corner buttons to open the variant track in table view. Open the table in split view by holding down the Ctrl key (Cmd on Mac) while clicking on the table icon (see the red arrow in figure 14). You will be able to see the variant track and the variant table in the same window.

The variant detection tool identifies around 8,000 variants in the bovine genome, about 13,000 in zebu, and about 23,000 in bison.

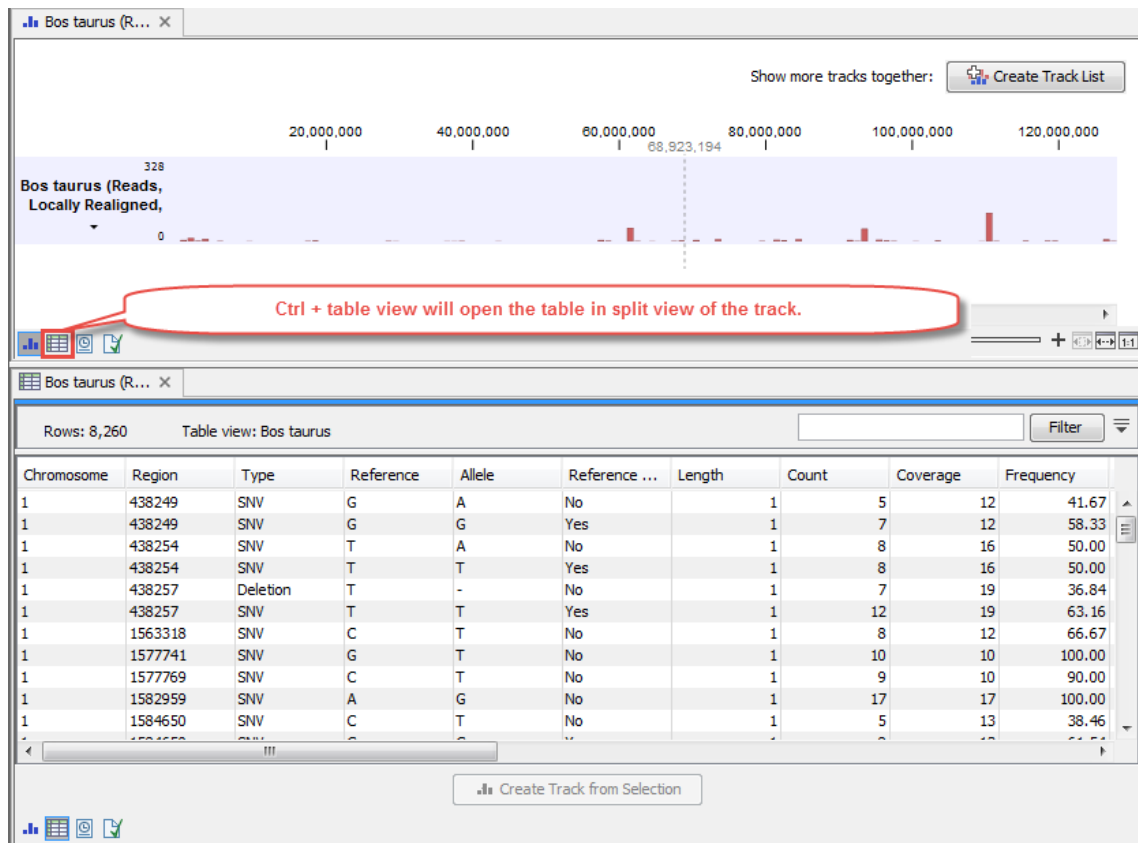




Figure 14: Open the variant table in split view.

Identify Species-specific Variants

We are now going to identify variants that are unique for the zebu species with the tool **Identify known mutations from sample mapping**.

1. Go to:

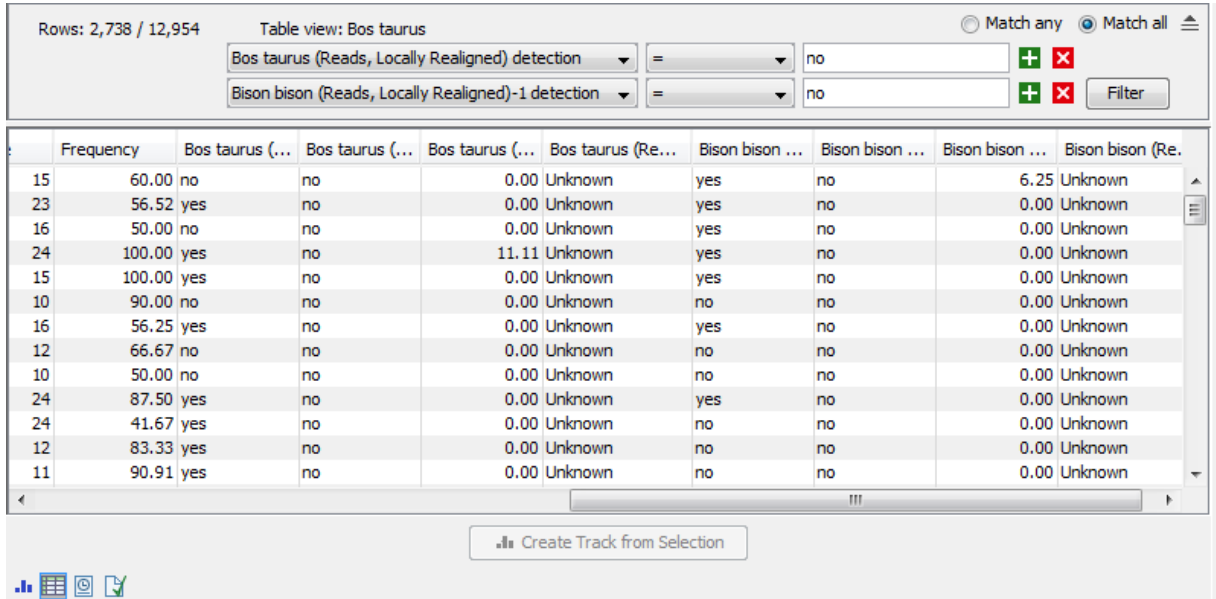
Toolbox | Resequencing Analysis  | **Identify known mutations from sample mapping** 

2. In the first wizard dialog, select the read mappings of *Bos taurus* and *Bison bison*. Click **Next**.

3. Select the zebu variant track **Bos indicus (Reads) - locally realigned (Variants)**. Leave the other parameters at their default value and click on the button labeled **Next**.

4. Create individual tracks and an overview track. Choose to Save them in a specified location (for example a new folder called zebu specific variants), and click on **Finish**.

5. Now open the file **Bos indicus (Reads, Locally Realigned, Variants) (MutationTest overview)**. We want to filter away from the table the variants that were also detected in the two other species, and keep only the ones that were not (the detection value is set to no). Open the filter and select "Bos taurus detection" in the first field, "=" in the second, and type "no" in the third field. Add an additional filter that you set up as you did the previous one but for "Bison bison detection" (see figure 15).



	Frequency	Bos taurus (...)	Bos taurus (...)	Bos taurus (...)	Bos taurus (Re...)	Bison bison ...	Bison bison ...	Bison bison ...	Bison bison (Re...
15	60.00	no	no	0.00	Unknown	yes	no	6.25	Unknown
23	56.52	yes	no	0.00	Unknown	yes	no	0.00	Unknown
16	50.00	no	no	0.00	Unknown	yes	no	0.00	Unknown
24	100.00	yes	no	11.11	Unknown	yes	no	0.00	Unknown
15	100.00	yes	no	0.00	Unknown	yes	no	0.00	Unknown
10	90.00	no	no	0.00	Unknown	no	no	0.00	Unknown
16	56.25	yes	no	0.00	Unknown	yes	no	0.00	Unknown
12	66.67	no	no	0.00	Unknown	no	no	0.00	Unknown
10	50.00	no	no	0.00	Unknown	no	no	0.00	Unknown
24	87.50	yes	no	0.00	Unknown	yes	no	0.00	Unknown
24	41.67	yes	no	0.00	Unknown	no	no	0.00	Unknown
12	83.33	yes	no	0.00	Unknown	no	no	0.00	Unknown
11	90.91	yes	no	0.00	Unknown	no	no	0.00	Unknown

Figure 15: You can filter for zebu specific variants in the overview table.

6. Click on the button "Filter", select all remaining rows and click on "Create Track from Selection". Save the new track in the Navigation Area.

The output from this analysis is a new variant track that holds only the variants (less than 3000) that are specific for the zebu species *Bos indicus* called "(Reads, Locally Realigned, Variants) (MutationTest overview)_selection". You can also see that you have different types of variation, some of which are heterozygous and others homozygous. You can export tables from the Workbench in Excel formats or text formats. Click on the table so that the table you would like to export is the selected view, then use the **Export** (📄) button in the toolbar to start the export. For small tables, you can also just **Copy** (📄) the contents of the table and paste into a spreadsheet for further processing.

We will now look only at the variants that are present in the coding part of the genome. We can do this with the tool **Filter Based on Overlap**.

1. Go to:

Toolbox | Track Tools (📊) | **Annotate and Filter** (🔍) | **Filter Based on Overlap** (🔗)

2. In the first wizard step, select the zebu specific variant track that you just saved as input. Click **Next**.
3. Select the **Bos taurus (Hereford)_CDS** track for overlap comparison and **Keep annotations that overlap**. Click **Next**.
4. Choose to save the results and click on the button labeled **Finish**.

The result is a variant track called "(Reads, Locally Realigned, Variants) (MutationTest overview)_selection (OF)" and containing a number in the range of 600 zebu-specific variants that occurs in the CDS regions.

Now we would now like to look only at variants that cause amino acid changes. We can do this with the tool **Amino Acid Changes**.

1. To run the analysis, go to the toolbox:

Toolbox | Resequencing Analysis (📁) | **Functional Consequences** (📁) | **Amino Acid Changes** (📁)

2. In the first wizard step, select the most recent zebu specific variant track (the one filtered based on overlap with CDS) as input. Click **Next**.
3. Select the CDS, mRNA, and sequence tracks as shown in figure 16 and choose to detect only non-synonymous changes by ticking the box "Filter synonymous". Leave also the option "Filter CDS regions with no variants" checked. Click **Next**.

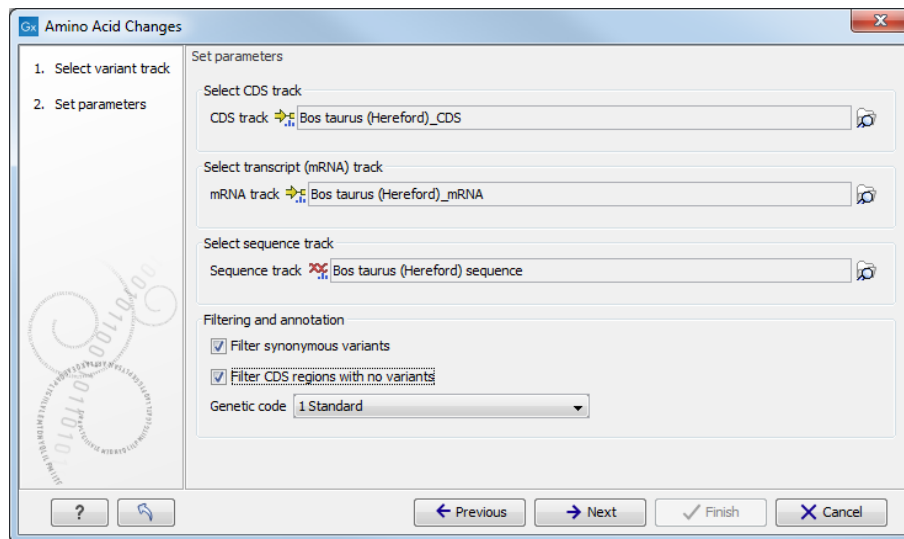


Figure 16: Select CDS, mRNA, and sequence tracks and check Filter synonymous and Filter CDS regions with no variants.

4. Choose to save the results and click on the button labeled **Finish**.

The result is a variant track called "(Reads, Locally Realigned, Variants) (MutationTest overview)_selection (OF, AAC)" with over 200 non-synonymous zebu-specific variants found in the CDS regions that cause amino acid changes.

Overrepresentation analysis

To gain more insight into the underlying biology we are now going to perform a gene ontology overrepresentation analysis with the tool **GO Enrichment Analysis**. The Gene Ontology site <http://www.geneontology.org/GO.downloads.annotations.shtml> contains the bovine gene ontology annotation.

1. Go to http://cvsweb.geneontology.org/cgi-bin/cvsweb.cgi/go/gene-association/gene_association.goa_cow.gz?rev=HEAD and download the bovine gene ontology annotation to your computer (1.6 MB).
2. Import the file into the Workbench with the import function found in the toolbar: **Import | Standard Import**

3. Select **Gene Ontology Annotation File** as **Format** in the import wizard (or alternatively use **Automatic Import**, which will give the same result). The file does not have to be unzipped before import.
4. To analyze if any biological pathways are overrepresented for the detected zebu-unique variants, go to the toolbox and run the tool:

Toolbox | Resequencing Analysis  | **Functional Consequences**  | **GO Enrichment Analysis** 

5. In the first wizard step, select the variant track that was created with the **Amino Acid Changes** tool **Bos indicus (Reads) - locally realigned (Variants, CTRL, OF, AAC)** as input. Click **Next**.
6. Select the parameter track **Bos taurus (Hereford)_Gene** and the GO annotation table, uncheck **Exclude computationally inferred GO terms**, and check **Allow gene name synonyms** as shown in figure 17. Click **Next**.

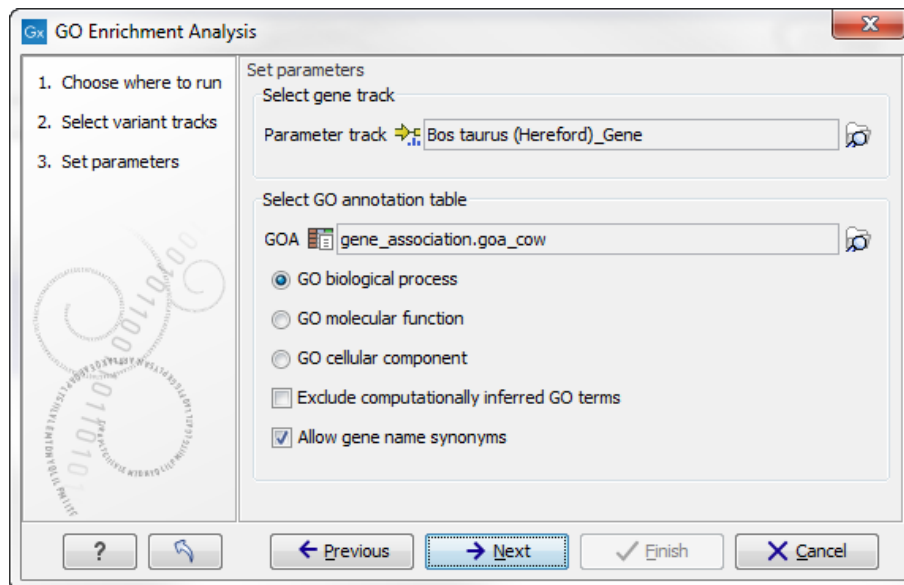


Figure 17: Select parameter track and GO annotation table as specified.

7. Choose to save the results and click on the button labeled **Finish**.

You get two different types of output: a **GO-enrichment table** and a variant track that is essentially the same variant table as the one generated in the **Amino Acid Changes** analysis but with five extra columns (see the red box in the side panel in figure 18).

Open the GO-enrichment table by clicking on the file name in the **Navigation Area**. We can sort the rows by the calculated p-value (from a hypergeometric distribution test comparing frequency of occurrence in a sample vs in all genes) by clicking once on the column header "P-values". It is clear that in the top rows the pathways related to the immune response dominate (figure 19). This matches previous findings showing that zebu responds differently to some infections and also has a better innate immune response compared to the taurine and the bison species [Freeman et al., 2008].

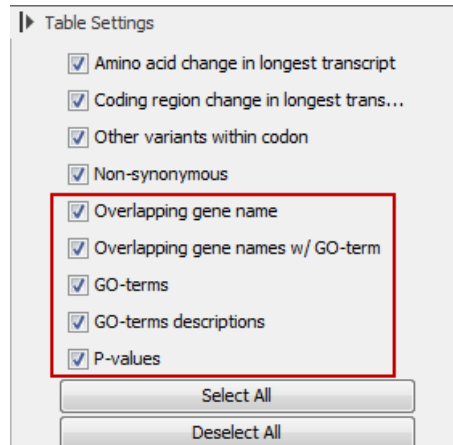


Figure 18: Five extra columns have been added compared to the Amino Acid Changes table.

GO-enrichment x

Rows: 7,994 GO enrichment analysis

GO term	Description	Occurrences in all genes	Occurrences in sample	P-values
0032757	positive regulation of interleukin-8 production	16	5	1.11E-5
0002755	MyD88-dependent toll-like receptor signaling pathway	9	4	1.80E-5
0034123	positive regulation of toll-like receptor signaling pathway	4	3	3.10E-5
0032755	positive regulation of interleukin-6 production	34	6	4.99E-5
0006955	immune response	137	11	9.16E-5
0042346	positive regulation of NF-kappaB import into nucleus	14	4	1.32E-4
0007252	I-kappaB phosphorylation	7	3	2.59E-4
0050707	regulation of cytokine secretion	7	3	2.59E-4
0042496	detection of diacyl bacterial lipopeptide	2	2	3.97E-4
0071726	cellular response to diacyl bacterial lipopeptide	2	2	3.97E-4
0002504	antigen processing and presentation of peptide or polysaccharide antigen via MHC class II	10	3	8.50E-4
0032722	positive regulation of chemokine production	11	3	1.15E-3
0032611	interleukin-1 beta production	3	2	1.18E-3
0045356	positive regulation of interferon-alpha biosynthetic process	3	2	1.18E-3
0071224	cellular response to peptidoglycan	3	2	1.18E-3
0050729	positive regulation of inflammatory response	25	4	1.41E-3
0044130	negative regulation of growth of symbiont in host	12	3	1.51E-3
0045087	innate immune response	89	7	1.98E-3
0001867	complement activation, lectin pathway	4	2	2.32E-3
0042773	ATP synthesis coupled electron transport	4	2	2.32E-3
0019221	cytokine-mediated signaling pathway	32	4	3.59E-3
0002024	diet induced thermogenesis	5	2	3.82E-3
0002830	positive regulation of type 2 immune response	5	2	3.82E-3
0042116	macrophage activation	5	2	3.82E-3
0008203	cholesterol metabolic process	35	4	4.99E-3
0002925	positive regulation of humoral immune response mediated by circulating immunoglobulin	6	2	5.65E-3
0045359	positive regulation of interferon-beta biosynthetic process	6	2	5.65E-3

Figure 19: Sort on p-value by clicking once on the column header.

Automating the analysis as workflows

Some of the steps in this analysis can be combined into workflows. This is shown in the two examples in (figure 20 and figure 21). These workflows, combined with batch mode of operation, greatly simplify the analysis, and can be used to repeat it, for instance, to look for similar variant enrichment in the other two species of bovines. For further information about how to create and run workflows, please see the tutorial "An Introduction to Workflows" that can be found here: <http://resources.qiagenbioinformatics.com/tutorials/Workflow-intro.pdf>.

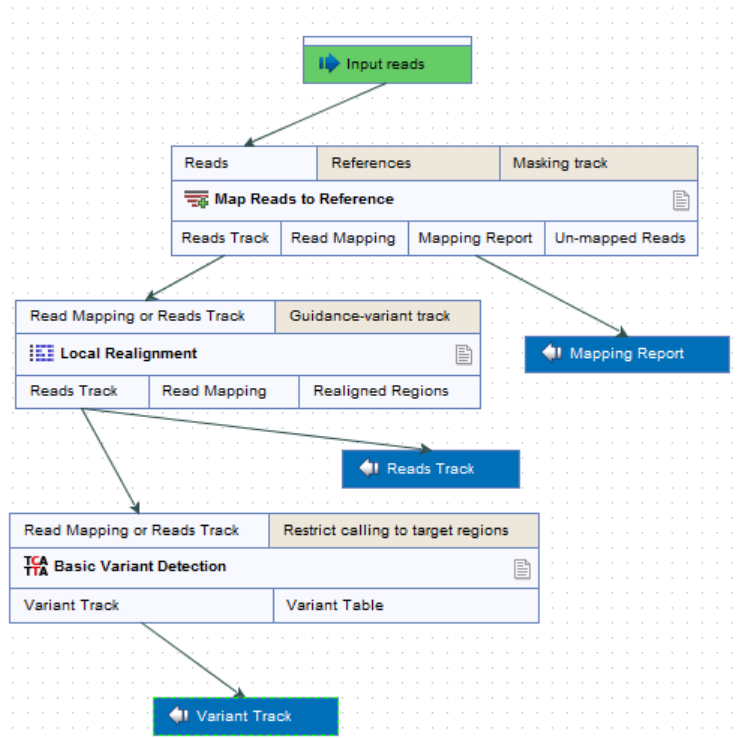


Figure 20: Mapping of the reads and variant calling workflow.

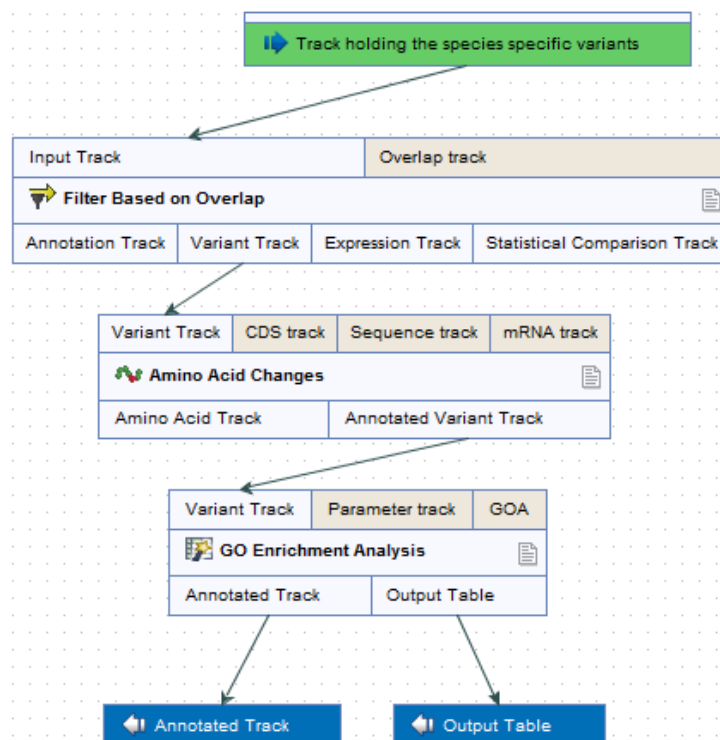


Figure 21: Variant filtering and annotation workflow.

Bibliography

- [Cosart et al., 2011] Cosart, T., Beja-Pereira, A., Chen, S., Ng, S. B., Shendure, J., and Luikart, G. (2011). Exome-wide dna capture and next generation sequencing in domestic and wild species. *BMC Genomics*, 12:347.
- [Freeman et al., 2008] Freeman, A. R., Lynn, D. J., Murray, C., and Bradley, D. G. (2008). Detecting the effects of selection at the population level in six bovine immune genes. *BMC Genet*, 9:62.