BIONUMERICS Tutorial:
# CFSAN SNP pipeline

## 1 Introduction

The ***CFSAN SNP pipeline*** is a SNP pipeline created by the FDA Center for Food Safety and Applied Nutrition (CFSAN) [1] which can be launched on sequence read set experiments in BIONUMERICS.

A ***CFSAN SNP pipeline*** analysis in BIONUMERICS consists of the following steps:

1. Choose an appropriate reference sequence for your samples of interest.

2. Create a comparison of your samples of interest.

3. Launch a wgSNP analysis with the ***CFSAN SNP pipeline*** on the Calculation Engine.

4. Import the results in BIONUMERICS and analyse the wgSNP clustering in the *Comparison* window.

The CFSAN SNP pipeline was developed for closely related organisms and is therefore not suited for the analysis of organisms for which no single appropriate reference sequence is available. For more details on the CFSAN SNP pipeline and the performed mapping, variant calling, SNP filtering and clustering steps, we refer to Davis *et al.* 2015 [1].

## 2 Preparing the database

The CFSAN SNP pipeline can only be performed in BIONUMERICS after installation of the *WGS tools plugin* in the BIONUMERICS database (***File*** > ***Install / remove plugins...*** ( )).

As the CFSAN SNP pipeline is only available on the Cloud Calculation engine make sure to select the options ***Use default Cloud Calculation Engine*** and ***Enable running jobs on Cloud Calculation Engine*** during installation of the *WGS tools plugin*. The Calculation engine option requires credits for running jobs on the Applied Maths cloud calculation engine. Credits are linked to credentials that you need to enter when installing the *WGS tools plugin*.
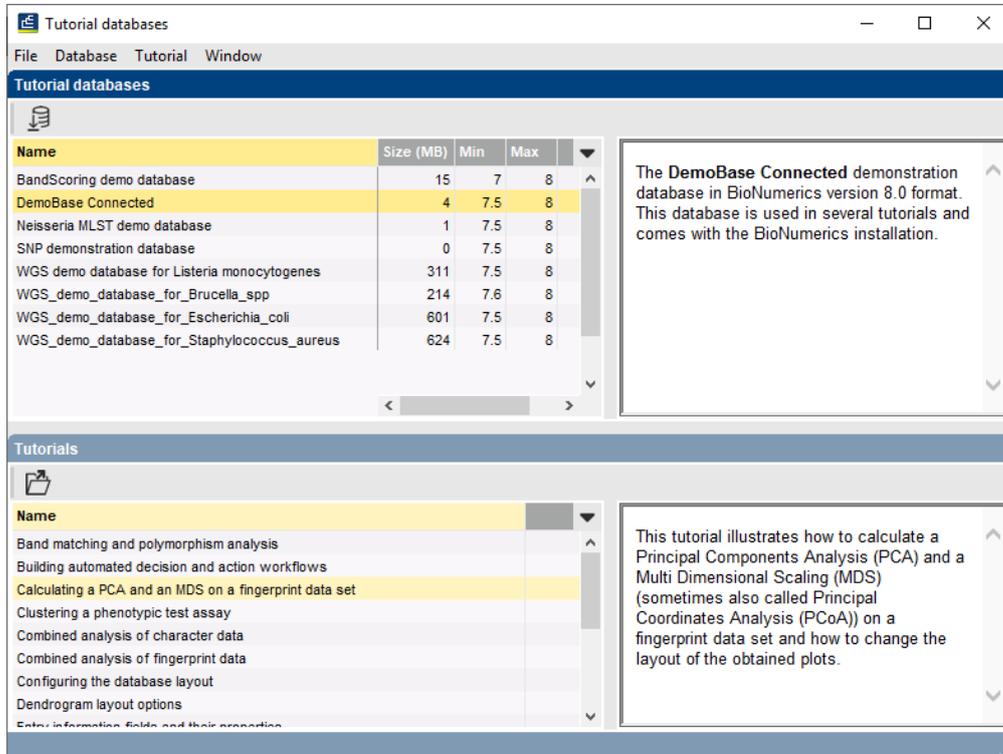
In this tutorial the **WGS demo database for *Salmonella*** will be used in which the *WGS tools plugin* is already installed. No credits are assigned to the demo project so no CFSAN SNP pipeline jobs can be launched on the external calculation engine. Please contact Applied Maths to obtain more information.

The **WGS demo database for *Salmonella*** can be downloaded directly from the *BIONUMERICS Startup* window (see 2.1), or restored from the back-up file available on our website (see 2.2)

## 2.1 Option 1: Download demo database from the Startup Screen

1. Click the ⬇ button, located in the toolbar in the *BIONUMERICS Startup* window.

This calls the *Tutorial databases* window (see Figure 1).



**Figure 1:** The *Tutorial databases* window, used to download the demonstration database.

2. Select **WGS_demo_database_for_Salmonella_enterica** from the list and select *Database* > *Download* ( 🗄 ).

3. Confirm the installation of the database and press <*OK*> after successful installation of the database.

4. Close the *Tutorial databases* window with *File* > *Exit*.

The **WGS_demo_database_for_Salmonella_enterica** appears in the *BIONUMERICS Startup* window.

5. Double-click the **WGS_demo_database_for_Salmonella_enterica** in the *BIONUMERICS Startup* window to open the database.

## 2.2 Option 2: Restore demo database from back-up file

A BIONUMERICS back-up file of the demo database for *Salmonella enterica* is also available on our website. This backup can be restored to a functional database in BIONUMERICS.

6. Download the file WGS_Salm.bnbk file from https://www.applied-maths.com/download/sample-data, under 'WGS_demo_database_for_Salmonella_enterica'.

✏️ In contrast to other browsers, some versions of Internet Explorer rename the `WGS_Salm.bnbk` database backup file into `WGS_Salm.zip`. If this happens, you should manually remove the `.zip` file extension and replace with `.bnbk`. A warning will appear ("If you change a file name extension, the file might become unusable."), but you can safely confirm this action. Keep in mind that Windows might not display the `.zip` file extension if the option "Hide extensions for known file types" is checked in your Windows folder options.

7. In the *BIONUMERICS Startup* window, press the 🗄️ button. From the menu that appears, select **Restore database...**.

8. Browse for the downloaded file and select **Create copy**. Note that, if **Overwrite** is selected, an existing database will be overwritten.

9. Specify a new name for this demonstration database, e.g. "WGS_Salmonella_demobase".

10. Click <**OK**> to start restoring the database from the backup file.

11. Once the process is complete, click <**Yes**> to open the database.

The *Main* window is displayed (see Figure 2).



**Figure 2:** The *Salmonella* demonstration database: the *Main* window.

# 3 Select a reference sequence

To demonstrate the CFSAN SNP pipeline analysis in BIONUMERICS a publicly available benchmark dataset (NCBI BioProject PRJNA412988) for which epidemiological data is available will be used. The sequence read set data of this BioProject was already imported in the demo database and used as input for a denovo assembly job and wgMLST analysis.

We will first import the closed genome sequence of *S. typhimurium* LT2 which is most referred to in literature and is also used in Saltykova *et al.* (2018) [2] for the same benchmark dataset. An appropriate reference sequence for wgSNP analysis should be selected with care and is preferably a high quality and complete genome sequence of an organism which is closely related to all samples in the dataset. It is recommended to remove all contigs which are smaller than 1000 bp from the reference sequence.

1. In the *Main* window, select **File** > **Import...** ( 📥 , **Ctrl+I**) to open the *Import* dialog box.

2. Choose the option **Download sequences from internet** under the **Sequence type data** item in the tree and click <**Import**> (see Figure 3).

3. Enter the accession code **NC_003197** in the **Accession codes** input field.

4. Choose one of the available download sites from the list, e.g. **NCBI**.

5. With the option **Preview sequences** checked, press <**Next**>.



**Figure 3:** Import sequences dialog box.

The import routine fetches the sequence from the selected database and shows detailed information in the next step.

6. Press <**Next**>.

The next step of the import wizard lists the templates that are present to import sequence information in the database. As this is the first time we import sequences from an online repository, we need to create a new import template by specifying **Import rules**.

7. Click <**Create new**> to create a new import template.

Each header tag (e.g. ID, AC, . . . ) corresponds to a row in the grid panel.

8. Select **AC - ACCESSION** in the list and click <**Edit destination**> or double-click on **AC - ACCESSION**. Select **Key**, and press <**OK**>.

9. Select **OS - SOURCE** in the list and click <**Edit destination**> or double-click on **OS - SOURCE**. Under **Entry info field** select **Organism**, and press <**OK**>.

The grid is updated (see Figure 4).

**Figure 4:** Import template.

10. Click <***Next***> and press <***Finish***>.

11. Specify a template name (e.g. **NCBI**) and optionally enter a description. Press <***OK***>.

12. Highlight the newly created template and select ***denovo*** as ***Experiment type*** (see Figure 5).
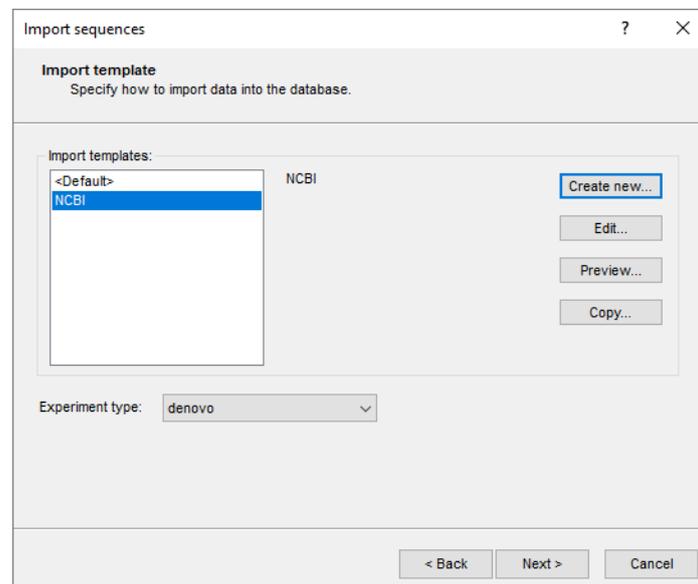


**Figure 5:** Import sequences dialog box.

13. Press <***Next***>.

The *Database links* wizard page will indicate that 1 new entry will be created during import.

14. Press <***Finish***>.

The sequence is imported in the database and is automatically selected.

15. Click on the green colored dot of the newly created entry in the *Experiment presence* panel to open the *Sequence editor* window.

The sequence is displayed in the upper panel and a graphical representation of the sequence is displayed in the panel below. The *Annotation* panel holds the GenBank features, and the header information is stored in the *Header* panel.

16. Close the *Sequence editor* window.
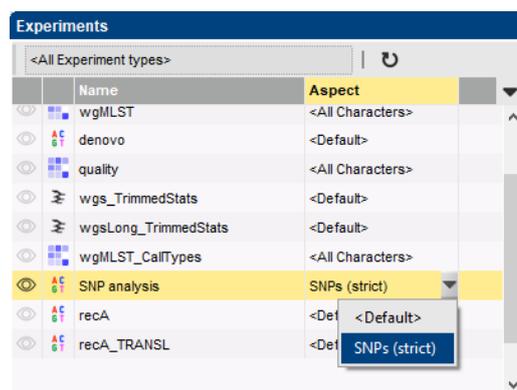
# 4 Create a comparison

We will now create a comparison containing the entries of the benchmark dataset (i.e. the samples of BioProject PRJNA412988).

1. In the *Database entries* panel of the *Main* window make sure no entries are selected. All entries can be deselected at once by pressing **Edit** > **Clear selection** (**Ctrl+Shift+A**).

2. Select **Edit** > **Find object in list...** (📇, **Ctrl+Shift+F**) and in the *Find* dialog box type "PR-JNA412988" and press <**Select all**>.

The 32 entries of the benchmark dataset are now selected.

3. Click on **Edit** > **Create new object...** ( + ) in the *Comparisons* panel to open the comparison window for the 32 selected entries.

A CFSAN SNP pipeline analysis can be launched on the calculation engine for the entries in the created comparison. Following the same principles as cluster analyses in comparisons, comparison jobs apply on the whole comparison and the data stored in the *active* experiment and (where applicable) the active aspect in the *Experiments* panel (see Figure 6 for an example).



**Figure 6:** The *Experiments* panel in the *Comparison* window. The experiment type highlighted in yellow (**SNP analysis** is the active experiment, the **SNPs (strict)** is the active aspect.

4. Save the comparison by selecting **File** > **Save** (💾, **Ctrl+S**). Specify a name for the comparison for example "CFSAN SNP analysis" and press <**OK**>.
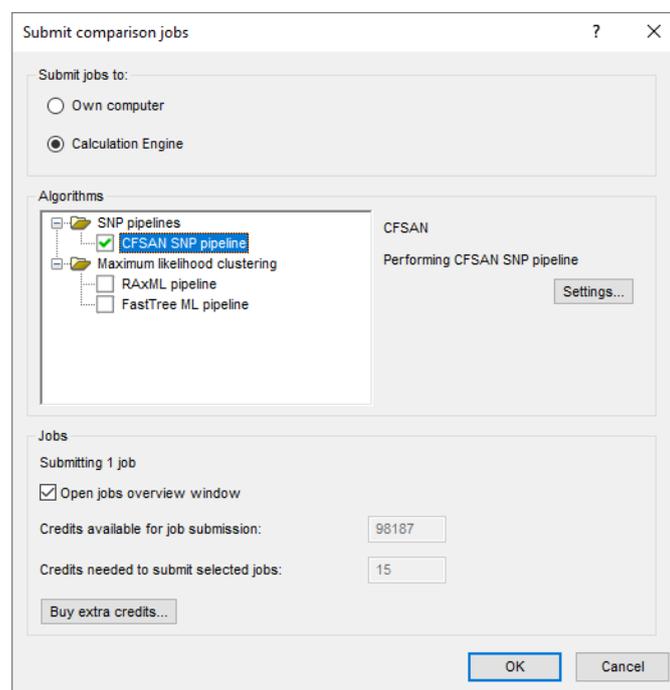
# 5 Launch the CFSAN SNP pipeline

1. In the *Experiments* panel select the **wgs** experiment which contains the sequence read set data of the entries in the *Comparison* window.

2. To launch the CFSAN SNP pipeline comparison job, select **File** > **Launch comparison jobs...** ( ▷ ).

In case the comparison was not saved to the database yet, you will be prompted to save first. Comparisons should be saved before any jobs can be launched.

Subsequently, the *Submit comparison jobs* dialog box will open (see Figure 7).

If the active experiment does not support any comparison jobs, an error message is shown and the *Submit comparison jobs* dialog box will not appear.



**Figure 7:** The *Submit comparison jobs* dialog box.

The CFSAN SNP pipeline is not available on the local calculation engine.

3. In the **Submit jobs to** panel select **Calculation engine** and in the **Algorithms** panel select the **CFSAN SNP pipeline** algorithm (see Figure 7).

4. With the **CFSAN SNP pipeline** algorithm highlighted, press the <**Settings...**> to open the *CFSAN SNP pipeline settings* dialog box in which the reference sequence(s) and other settings for the CFSAN SNP pipeline job can be defined (see Figure 8).

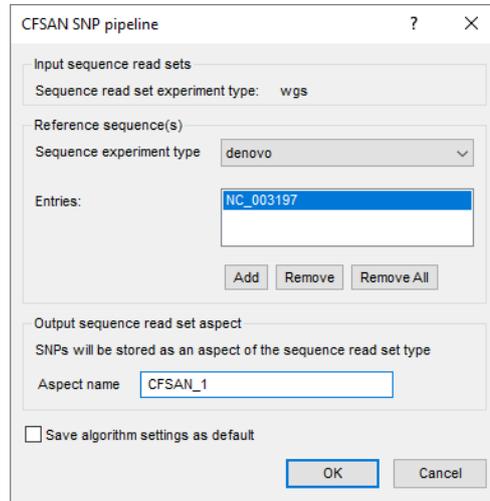The **Sequence read set experiment type** selected as **Input sequence reads** is read-only because it was set by choosing the active experiment prior to calling the *Submit comparison jobs* dialog box.

Under **Reference sequence(s)**, the reference genome that should be used for the mapping needs to be specified.

5. Select the **denovo** sequence experiment type from the sequence experiment type drop-down list.

6. Press the <***Add***> button and specify the entry that contains the reference genome i.e. the entry with Key NC_003197 (see Figure 8).

If needed, multiple reference genomes can be chosen by repeating this step.



**Figure 8:** The *CFSAN SNP pipeline settings* dialog box.

The resulting SNP matrix (i.e. the output from the CFSAN SNP pipeline) will be stored as an aspect of the input sequence read set experiment type. An ***Aspect name*** can be entered manually or the default name can be accepted.

To avoid having to re-enter the above settings for a subsequent analysis, one can save them as defaults to the database with ***Save algorithm settings as default***.

7. Select <***OK***> to confirm the settings and close the *CFSAN SNP pipeline settings* dialog box.

8. In the *Submit comparison jobs* dialog box leave the option ***Open jobs overview window*** selected and select <***OK***> to launch the CFSAN SNP pipeline job.

The job is submitted to the Calculation Engine and the *Job overview* window opens. In the *Job overview* window, the job type, job name, time of submission, job status, a description of the job, its progress and much more can be monitored.
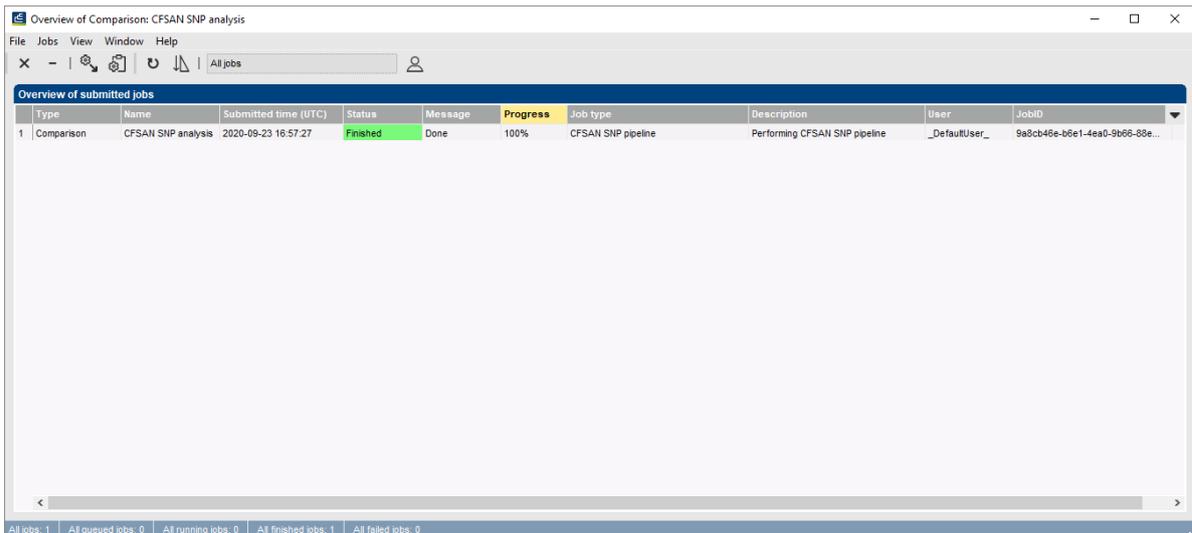
# 6 Import and analyse the CFSAN SNP pipeline results

Once the job has been finished (see Figure 9), the results can be imported in the database by selecting ***Jobs*** > ***Get results*** (⚙) from the *Job overview* window.
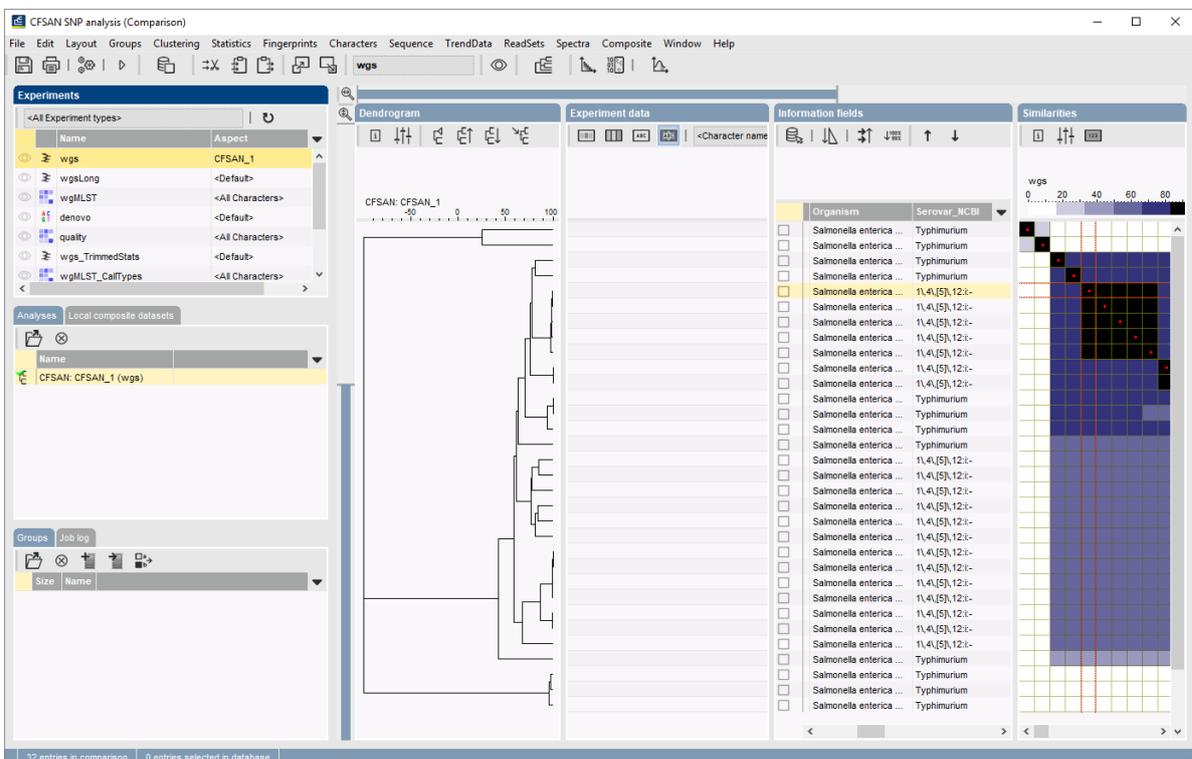
1. When the job is finished, highlight the job and select ***Jobs*** > ***Get results*** (⚙) to import the results in the comparison window.

2. Close the *Job overview* window.

The *Comparison* window now looks like Figure 10.

In the *Experiments* panel, the aspect containing the CFSAN SNP pipeline results is selected. The default name of the first CFSAN aspect is "CFSAN_1"). The *Analyses* panel lists the performed analyses (for a CFSAN analysis, the analysis type (CFSAN), the aspect (CFSAN_1) and the ex-

**Figure 9:** The *Job overview* window listing a finished CFSAN SNP pipeline job.



**Figure 10:** The *Comparison* window after import of the CFSAN SNP pipeline job results.

periment type on which the analysis was performed (wgs)is included in the name of the analysis).

3. Click on the 👁 next to the experiment name **wgs** in the *Experiments* panel to display the SNP matrix in the *Experiment data* panel.

The SNP matrix is now visible in the *Experiment data* panel. By default, the character name is present above the SNP positions but this can be changed to the reference sequence ID, the contig number, the position on the contig or the position on the reference sequence (see Figure 11).

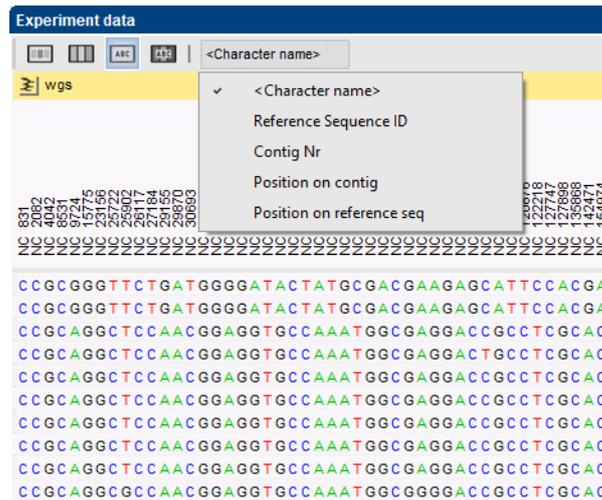The default dendrogram display settings in BIONUMERICS are set to similarity values. However,

**Figure 11:** The *Experiment data* panel.

the CFSAN SNP pipeline calculates pairwise SNP distances for clustering. The dendrogram display settings should therefore be adjusted.

4. Select **Clustering** > **Dendrogram display settings...** ( ) to open the *Dendrogram display settings* dialog box.

5. Select **Use distances** and **Show branch distances** (see Figure 12). Press <**OK**>.
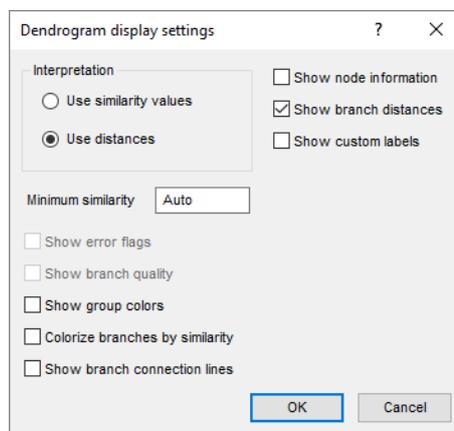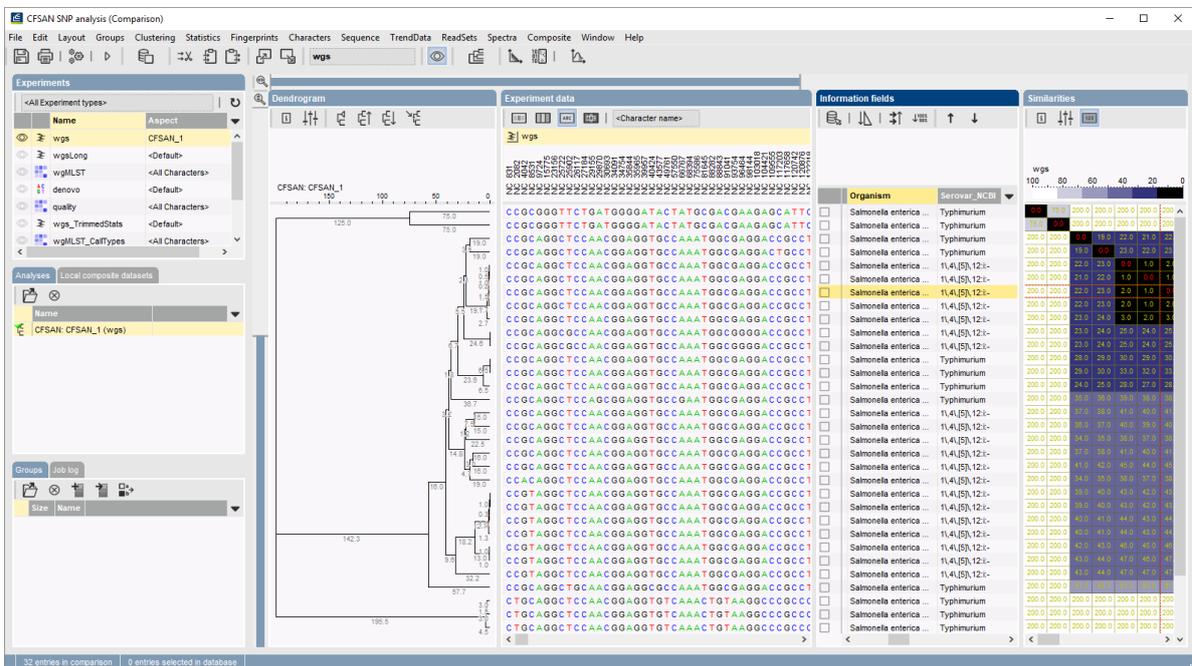


**Figure 12:** The dendrogram display settings.

6. Select **Clustering** > **Similarity matrix** > **Show values** ( ) to visualise the pairwise SNP distances in the *Similarities* panel.

The *Comparison* window now looks like Figure 13.

The log file of the CFSAN SNP pipeline job can be consulted in the Job log panel (see Figure 14).

**Figure 13:** The *Comparison* window after visualising the SNP matrix, adjusting the dendrogram display settings and visualising the pairwise SNP distances.



**Figure 14:** The Job log panel.

# Bibliography

[1] Steve Davis, James B Pettengill, Yan Luo, Justin Payne, Al Shpuntoff, Hugh Rand, and Errol Strain. Cfsan snp pipeline: an automated method for constructing snp matrices from next-generation sequence data. *PeerJ Computer Science*, 1:e20, 2015.

[2] A. Saltykova, V. Wuyts, W. Mattheus, S. Bertrand, NHC Roosens, and K. Marchal. Comparison of snp-based subtyping workflows for bacterial isolates using wgs data, applied to salmonella enterica serotype typhimurium and serotype 1,4,[5],12:i:-. *PLoS ONE*, 13(2), 2018.