BIONUMERICS Tutorial:

# wgSNP with mapping performed locally on your own computer

## 1 Introduction

A Single Nucleotide Polymorphism (SNP) is a variation in a single nucleotide, which occurs at a specific position of the genome. When performed on whole genome sequences (WGS), this analysis is referred as **whole genome SNP (wgSNP) analysis**. When performed in BIONUMERICS a typical workflow for wgSNP analysis consists of following steps:

1. Choose a reference sequence

2. Map sequence reads against the reference sequence (locally or on the cloud calculation engine)

3. Perform wgSNP analysis and filter out relevant SNPs

4. wgSNP clustering

In this tutorial we will focus on the first two steps of the workflow, with local mapping of the reads against the reference sequence. The last two steps are covered in the tutorial "wgSNP: analysis and clustering".

## 2 Preparing the database

Mapping of reads against the reference sequence can only be performed locally on your own computer after installation of the *WGS tools plugin* in the BIONUMERICS database (**File** > **Install / remove plugins...** ( )).

During installation of the *WGS tools plugin* plugin, make sure to select the options **Use default Cloud Calculation Engine** and **Local calculations only**.

The **WGS demo database** for *Staphylococcus aureus* already contains the installed *WGS tools plugin*. This demo database 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. To download the database directly from the *BIONUMERICS Startup* window, click the ⬇ button, located in the toolbar in the *BIONUMERICS Startup* window.

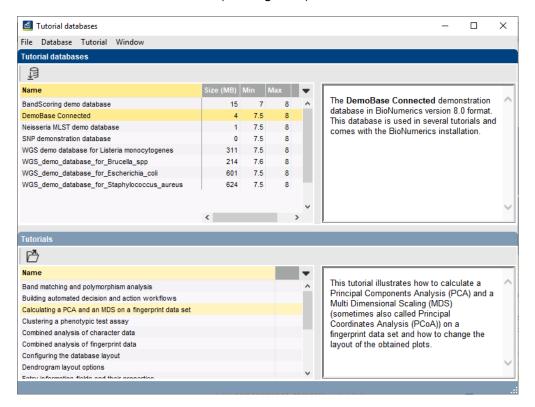This calls the *Tutorial databases* window (see Figure 1).



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

2. Select the **WGS_demo_database_for_Staphylococcus_aureus** 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_Staphylococcus_aureus** appears in the *BIONUMERICS Startup* window.

5. Double-click the **WGS_demo_database_for_Staphylococcus_aureus** 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 whole genome demo database for Staphylococcus aureus is also available on our website. This backup can be restored to a functional database in BIONU-MERICS.

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

In contrast to other browsers, some versions of Internet Explorer rename the `wgMLST_SAUR.bnbk` database backup file into `wgMLST_SAUR.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 [icon] 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. "Whole genome Staphylococcus aureus demobase".

10. Click <***OK***> to start restoring the database from the backup file (see Figure 2).
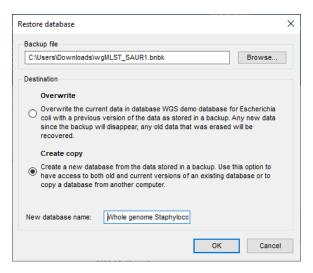


**Figure 2:** Restoring the whole genome demonstration database from the BioNumerics backup file `wgMLST_SAUR.bnbk`.

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

The *Main* window is displayed (see Figure 3).

## 3    Creating a reference mapped sequence type

First, we will create a sequence type to store the reference mapped sequences in:

1. Click on the *Experiment types* panel to activate it and select ***Edit*** > ***Create new object...*** ( + ).

2. From the *Create a new experiment type* dialog box that pops up, select ***Sequence type*** and press <***OK***>.

3. In the *New sequence type* wizard, enter a ***Sequence type name*** (e.g. "My wgSNP").

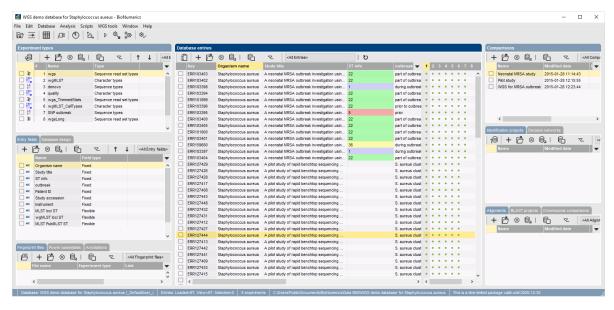4. Press <***Next***> (see Figure 4).

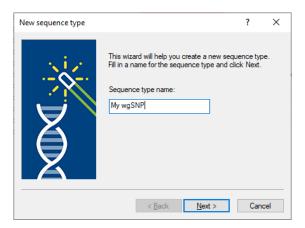**Figure 3:** The *Staphylococcus aureus* demonstration database: the *Main* window.



**Figure 4:** New sequence type experiment.

5. Leave the default **Nucleic acid sequences** option checked and check the option **Use reference sequence as mapping template** (see Figure 5).
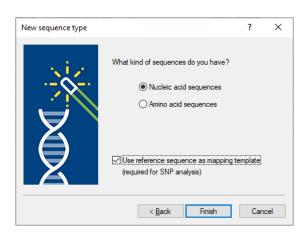


**Figure 5:** Use as mapping template.

6. Press <***Finish***> to create the reference mapped sequence type.

The reference mapped sequence type is added to the *Experiment types* panel.

# 4   Assigning a reference sequence

The first sequence that is imported in the reference mapped sequence type will automatically be assigned as the reference. The reference sequence might be a closed, fully annotated genome sequence (e.g. downloaded from an online repository such as NCBI), but could as well be a de novo assembled sequence, consisting of multiple contigs (i.e. a draft genome).

> Tutorials are available on our website explaining how to import genome sequences (***File*** > ***Import...*** ( 📷 , **Ctrl+I**), "Sequence data") and how to import and assemble sequence reads (***File*** > ***Import...*** ( 📷 , **Ctrl+I**), "Sequence read sets").

In this tutorial, we will copy the sequence from the ***denovo*** experiment type of entry ***ERR103401*** to our new sequence type experiment so it serves as a reference sequence:

1. Click on the green dot in the *Experiment presence* panel that corresponds to the ***denovo*** experiment for entry ***ERR103401***. In default configuration, the ***denovo*** experiment corresponds to the third column in the *Experiment presence* panel.

2. In the *Sequence editor* window that opens, select ***File*** > ***Save as...*** (**Ctrl+Shift+S**).

3. Highlight ***My wgSNP*** in the list and press <***OK***> (see Figure 6).
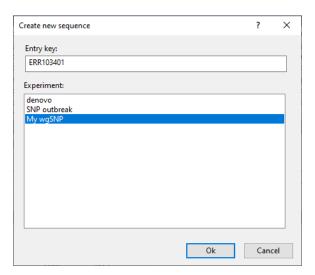


**Figure 6:** Save sequence to a new sequence type experiment.

4. Close the *Sequence editor* window.

We can check if this sequence is indeed used as reference:

5. In the *Experiment types* panel, double-click on ***My wgSNP*** to open its *Sequence type* window: the reference sequence is displayed here (see Figure 7).
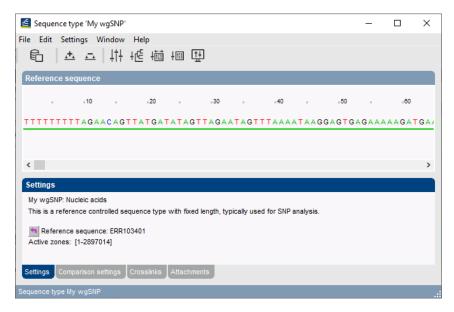
6. Close the *Sequence type* window.

**Figure 7:** The new sequence type experiment and the reference sequence.

# 5 Importing sequence read sets

Entries that we want to analyze using wgSNP need a sequence experiment that is obtained via mapping to the reference sequence, as this ensures collinearity of the sequences. For this, we need to start from the corresponding sequence read sets.

The preferred way of importing sequence read sets in BIONUMERICS is with the ***Import sequence read set files as links*** import routine (***File > Import...*** (⬚, **Ctrl+I**)). Importing sequence read sets as links is only possible when the *WGS tools plugin* is installed in the BIONUMERICS database (***File > Install / remove plugins...*** (⬚)).

In our demonstration database, the *WGS tools plugin* is already installed and all read sets are imported as links:

1. Click on the green colored dot for one of the entries in the first column in the *Experiment presence* panel. Column 1 corresponds to the first experiment type listed in the *Experiment types* panel, which is **wgs** in the default configuration.

The data link is displayed in the *Sequence read set experiment* window (see Figure 8).

2. Close the *Sequence read set experiment* window.

> ✎ Sequence read sets can also be stored inside the database with the ***Import sequence read set files*** import routine in the Import tree (***File > Import...*** (⬚, **Ctrl+I**)). One disadvantage of this option is the storage size of the sequence read sets.

# 6 Mapping sequence reads against the reference sequence

1. Select a few entries in the *Database entries* panel to include in the SNP analysis, using the ballot boxes next to the entries. Selected entries are marked by a checked ballot box: ☑.

2. Select ***WGS tools > Submit jobs...*** ( ▷ ) to display the *Submit jobs* dialog box (see Figure 9).
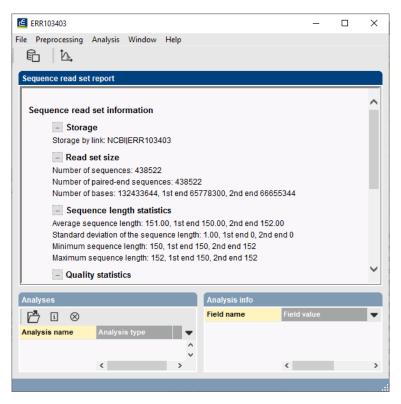
**Figure 8:** Data link.

The **Submit jobs to** option deals with the location where the jobs will be sent to for execution: either to your **Own computer** or to the cloud **Calculation Engine**.

     3. Since this tutorial covers wgSNP with mapping performed on the local computer, make sure the **Own computer** option is selected.

     4. Select the **My wgSNP** option as **Reference mapping** option (see Figure 9).

     5. With the **My wgSNP** option highlighted, press the <**Settings**> button.

The only reference mapper available is **SNAP**.

     6. Close the *Reference mapping settings* dialog box.

Since no credits are required to run jobs on your own computer, this section in the *Submit jobs* dialog box can be ignored.

     7. Press <**OK**> to launch the jobs.

When a reference mapping is already present for the submitted entries (and the **Re-submit already processed data** option was unchecked in the *Submit jobs* dialog box), an information message will appear, saying that no jobs are submitted.

The *Job overview* window will open, in which the status of the jobs can be followed (see Figure 11).

There are two options available in the *Job overview* window to import the job results in your BIONUMERICS database:

     8. Finished jobs can be imported with a manual action (**Jobs** > **Get results** (🔍)) or through an automatic update: select **File** > **Settings**, check both options and specify an interval (e.g. 10 min).
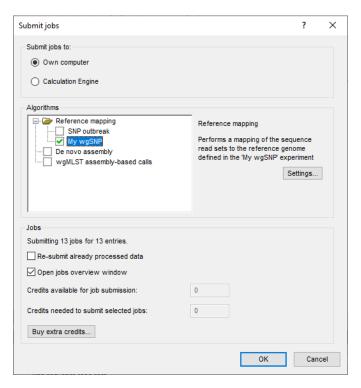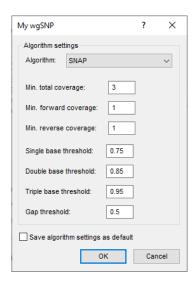
**Figure 9:** The *Submit jobs* dialog box.



**Figure 10:** Reference mapper.

9. Close the *Job overview* window.

The results of the mapping, i.e. the consensus sequences for the samples in the same frame as the reference sequence, are stored in the mapped sequence experiment type (here: **My wgSNP**).
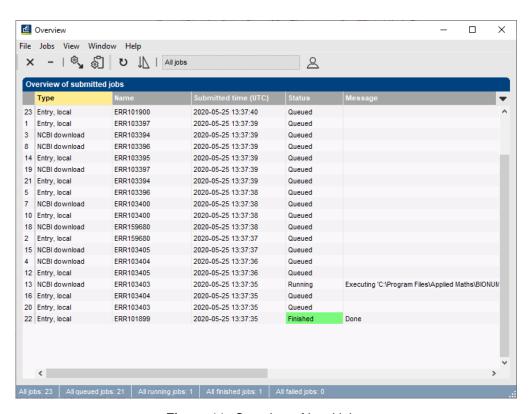
**Figure 11:** Overview of local jobs.