BIONUMERICS Tutorial:

# Fast character-based identification

## 1    Introduction

The BIONUMERICS software offers a tool for screening an entry against the database, based upon a character type experiment. This identification tool benefits from a bulk-fetching mechanism, which makes it many times faster for identification against large databases.

This tutorial describes how to perform such character-based identification of entries.

## 2    Preparing the database

### 2.1    Introduction to the demonstration database

We provide a **WGS demo database** for *Listeria monocytogenes* containing sequence read set data links for 51 samples, calculated de novo assemblies and wgMLST results (allele calls and quality information).

> The *de novo* assembly, wgMLST workflow and results will not be discussed in this tutorial.

The **WGS demo database** for *Listeria monocytogenes* can be downloaded directly from the *BIONUMERICS Startup* window (see 2.2), or restored from the back-up file available on our website (see 2.3).

### 2.2    Option 1: Download the 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).

2. Select the **WGS demo database for Listeria monocytogenes** 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*.

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

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

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

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

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

6. Download the file `WGS_LM01.bnbk` file from https://www.bionumerics.com/download/sample-data, under 'WGS_demo_database_for_Listeria_monocytogenes'.

> In contrast to other browsers, some versions of Internet Explorer rename the `WGS_LM01.bnbk` database backup file into `WGS_LM01.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 and make sure the name does not contain any spaces to ensure the successful installation of the *Listeria functional genotyping plugin*. Specify for example: "WGS_Listeria_demobase".

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



**Figure 2:** Restoring the **WGS demonstration database** from the backup file `WGS_LM01.bnbk`.

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

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



**Figure 3:** The *Listeria monocytogenes* demonstration database: the *Main* window.

# 3  About the demonstration database

The WGS demo database contains links to sequence read set data on NCBI's sequence read archive (SRA) for 51 publicly available sequencing runs. Additional information, stored in entry info fields (CollectionDate, CollectedBy, serovar, etc.) was collected from the corresponding publications and added to the demo database.

Seven experiments are present in the demo database and are listed in the *Experiment types* panel (see Figure 4).



**Figure 4:** The *Experiment types* panel in the *Main* window.

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 *Sequence read set experiment* window, the link to the sequence read set data on NCBI (SRA) with a summary of the characteristics of the sequence read set is displayed: *Read set size*, *Sequence length statistics*, *Quality statistics*, *Base statistics* (see Figure 5).

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

3. Click on the green colored dot for one of the entries in the third column in the *Experiment presence* panel. Column 3 corresponds to the third experiment type listed in the *Experiment types* panel, which is **denovo**.

The *Sequence editor* window opens, containing the results from the de novo assembly algorithm, i.e. concatenated de novo contig sequences (see Figure 6).

4. Close the *Sequence editor* window.

The sequence read set experiment type **wgs_TrimmedStats** contains some data statistics about the reads retained after trimming, used for the de novo assembly.

The sequence read set experiment type **wgsLong** contains the links to long read sequence read data (typically PacBio or MinION datasets). In this demo database, no links are defined for this experiment.

The other three experiments contain data related to the wgMLST analysis performed on the samples:

• Character experiment type **wgMLST** contains the allele calls for detected loci in each sample, where the consensus from assembly-based and assembly-free calling resulted in a single allele ID.

• Character experiment type **quality** contains quality statistics for the raw data, the de novo assembly and the different allele identification algorithms.

• Character experiment type **wgMLST_CallTypes**: contains details on the call types.

**Figure 5:** The sequence read set experiment card for an entry.

# 4 Fast identification based on characters

## 4.1 General identification

We will use the allele calls of the wgMLST profile stored in the **wgMLST** character experiment of entry 'SRR1695834' to screen for similar wgMLST profiles in the database.

1. In the *Main* window, click on ***Edit*** > ***Find object in list...*** (, **Ctrl+Shift+F**).

2. Type 'SRR1695834' (without the quotes) in the text field of the *Find* dialog box and click <***Select all***>.

The entry 'SRR1695834' is now selected in the database.

3. Click on the green dot corresponding with the **wgMLST** character experiment for the selected entry 'SRR1695834'.

The character experiment card, containing the different allele calls for all the wgMLST loci will open.

4. Close the experiment card by clicking on the triangle in the upper left corner of the card.

**Figure 6:** The *Sequence editor* window.

5. Choose **Analysis** > **Fast matching** > **Fast character matching...** in the *Main* window.

The *Fast character set matching* dialog box will open (see Figure 7).



**Figure 7:** The *Fast character set matching* dialog box.

6. Change the **Distance type** to **Categorical**.

7. Leave the other parameters at their default values and click <**OK**>.

More information on the other parameters can be found in the reference manual.

The *Fast matching* window will open showing the results of the screening (see Figure 8). In the *Entries* panel you will see the entries and their corresponding information fields that were used to

screen the database. In this case, we only used one entry (entry 'SRR1695834') for the screening. If you would use multiple entries, all these entries will be shown here.



**Figure 8:** The *Fast character set matching* dialog box.

8. Click on the entry 'SRR1695834' to show its matches.

In the *Matches* panel of the *Fast matching* window the matches for the selected entry are shown. In this case 5 matches should be present, one of which is a 100% match with entry 'SRR1695834' which we used for the screening (see Figure 8). The column **Distance** indicates the difference in distance between your entry and the match. In this particular case, this number can be interpreted as the number of allelic differences between the wgMLST profiles.

By looking at the values in the **CollectedBy** column we can see that matches with strains collected by both the CDC and FDA were detected (see Figure 8). In the next part of this tutorial you will learn to limit the search for matches to a certain subset of entries that have a certain (combination of) information field value(s) (see 4.2).

The results of the *Fast matching* window can also be exported to a comma separated text file.

9. Choose **File** > **Export...** ( 🖫 ) to export the results to text format.

10. Close the *Fast matching* window.

## 4.2 Limiting identification to a certain subset of entries in the database

This part of the tutorial will describe how you can limit the search for matches to a certain subset of entries that have a certain (combination of) information field value(s) by using a SQL statement. We will continue working on the previous example and try to obtain only matches of strains which were collected by the FDA.

> The SQL query will work only if you use the information field ID instead of the field name.

The information field IDs can be obtained from the *Entry fields* panel in the *Main* window window.

11. Click on the downward arrow to the right of the table headers in the *Entry fields* panel.

12. From the menu, select *'Set active fields...'* or press **Ctrl+F5**.

The *Set active fields* dialog box dialog will open (see Figure 9).



**Figure 9:** The *Set active fields* dialog box.

13. Check the box in front of **ID** (see Figure 9) and press <**OK**> to close the *Set active fields* dialog box.

The information field ID will now be visible in the *Entry fields* panel (see Figure 10). The ID of the **CollectedBy** information field is **COLLECTEDBY**.



**Figure 10:** The 'ID' field is now active in the *Entry fields* panel.

14. Make sure entry 'SRR1695834' is still selected in the *Main* window.

15. Choose **Analysis** > **Fast matching** > **Fast character matching...** in the *Main* window.

The *Fast character set matching* dialog box will open (see Figure 7).

16. Change the **Distance type** to **Categorical**.

17. In the **Database query** field, type the SQL statement "COLLECTEDBY"='FDA' (see Figure 11).

Two SQL statements can be combined by 'AND', when both statements should be fulfilled, or 'OR', when either statement should be fulfilled.

**Figure 11:** The SQL statement in the *Fast character set matching* dialog box.

18. Press <**OK**> to execute the fast matching.

The *Fast matching* window will open showing the results of the screening (see Figure 12). You will notice that now only 2 matches were found, both of strains which were collected by the FDA.

19. Close the *Fast matching* window.

**Figure 12:** The expected results in the *Fast matching* window.