BIONUMERICS Tutorial:

# Minhashing based cluster analysis of sequences

## 1 Aim

Similarity and distance-based trees (e.g. UPGMA and Neighbor joining) and phylogenetic trees (e.g. Maximum likelihood and Maximum parsimony trees) can be calculated in the *Comparison* window in BIONUMERICS, reflecting the relationships between the analyzed sequences. This tutorial covers clustering of sequences using a minhashing based calculated similarity.

MinHash techniques allow the comparison of large datasets of genomic sequences which is currently infeasible with alignment based approaches. During minhashing the k-mer content of each sequence is first determined and each k-mer is then passed through a hash function to obtain hashes. Retaining the lowest hashes enables the sampling of a random set of k-mers, which is called a "sketch" or "MinHash signature". Using only these signatures the similarity of the original sequences can be compared in a rapid and accurate way.

## 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 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).

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

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**.

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 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] 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_LMO1.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    Performing a minhashing based clustering of sequences

A minhashing based clustering is performed in the *Comparison* window.

1. In the *Database entries* panel of the *Main* window, select all entries that have an associated **denovo** experiment. To select all entries at once, use the **CTRL+A** shortcut combination.

2. Highlight the *Comparisons* panel in the *Main* window and select **Edit** > **Create new object...** ( + ) to create a new comparison for the selected entries.

3. Click on the experiment name **denovo** in the *Experiments* panel to highlight it.

✎    Clicking on the 👁 next to the experiment name **denovo** in the *Experiments* panel will display the sequences in the *Experiment data* panel. However, when analysing a lot of sequences at the same time, it could take quite a while before all sequences are displayed.

There is no need to align the sequences first. The similarities between all the sequences are calculated from the sequences' k-mers by the minhashing algorithm.

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

4. Select ***Clustering*** > ***Calculate*** > ***Cluster analysis (similarity matrix)...*** or press 🔲 and select ***Calculate cluster analysis***.

5. Choose ***Minhashing based*** with the default settings and press <***Next***>.

The **K-mer size** allows to specify a particular k-mer size for the minHashing algorithm to use. It is advised to use k = 31 or k = 51. A k-mer size with k = 31 is more sensitive at the genus level and k = 51 is more stringent, leading to fewer false positive matches. The amount of hashes that will be calculated can either be specified directly or can depended on the specified sourmash scaling factor, depending on whether the **Use sourmash scaling factor** option is checked or not.

✏️ More information on the different minHashing options can be found in the reference manual.

6. In the second step of the wizard, select ***Neighbor joining*** as dendrogram type and click <***Finish***>.

After the minHashing has been performed, the dendrogram and similarity matrix will be shown in the *Comparison* window (see Figure 7).

7. Click on the 🔳 in the *Similarities* panel to show the similarity values between the sequences.

8. Click on the 🔳 in the *Dendrogram* panel to show the settings that have been used for calculation of the minHashing based similarities and the dendrogram.

9. Save the comparison with ***File*** > ***Save*** (🖫, **Ctrl+S**).

**Figure 7:** Neighbor joining based on minhashing.