BIONUMERICS Tutorial:

# Annotating sequences by Prokka

## 1 Introduction

This tutorial demonstrates the annotation of genomes by Prokka [1] in BIONUMERICS. The Prokka algorithm is only available on the Calculation Engine.

To avoid issues with gaps in reference mapped (i.e. aligned) sequences, the Prokka algorithm is limited to sequence types that are *not* reference mapped.

## 2 Preparing the database

### 2.1 Introduction

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

As the Prokka job 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 Prokka 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 2.2), or restored from the back-up file available on our website (see 2.3).

### 2.2 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).

2. Select **WGS_demo_database_for_Salmonella_enterica** from the list and select **Database** > **Download** ( ).
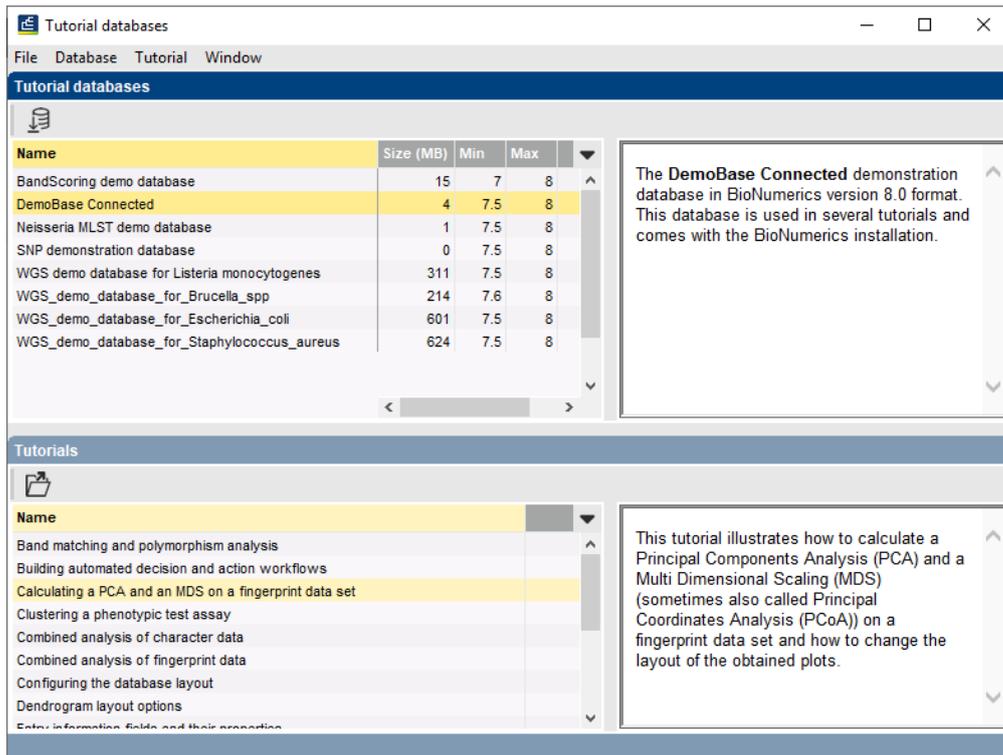
**Figure 1:** The *Tutorial databases* window, used to download the demonstration database.

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.

## 2.3   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.
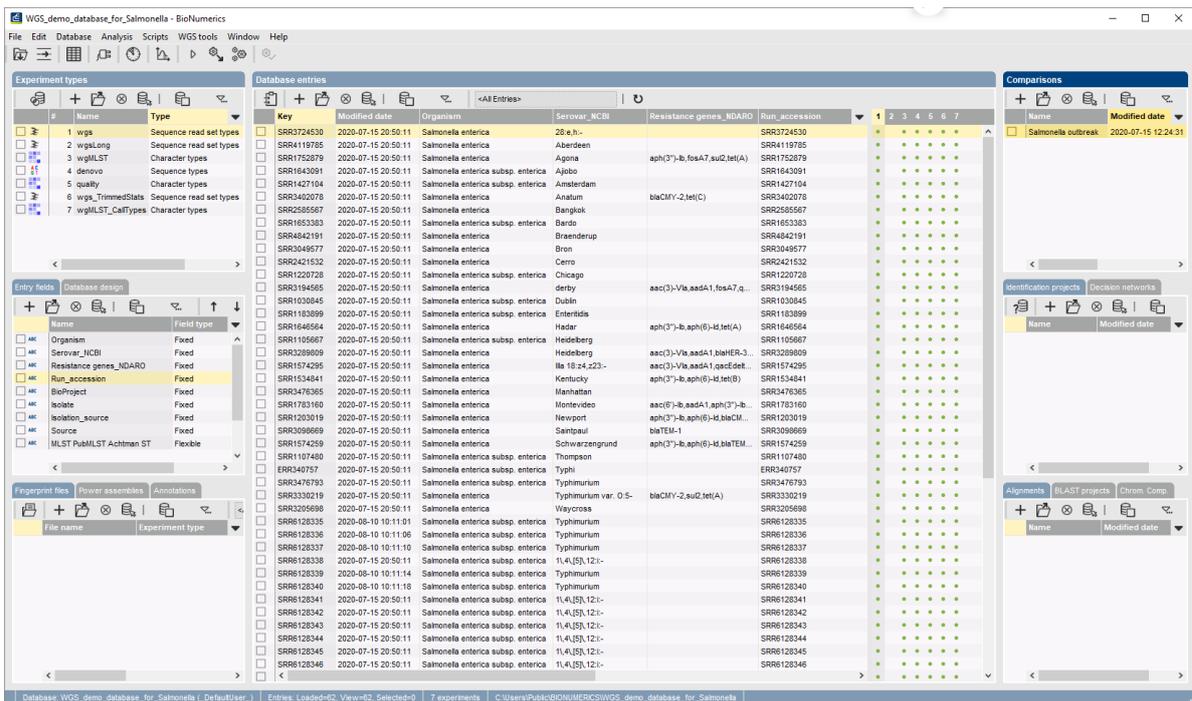
The *Main* window is displayed (see Figure 2).



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

# 3 Launch a Prokka job on the Calculation Engine

The **Annotation by Prokka** pipeline can be launched from the *Main* window on one or multiple selected entries.

1. In the *Main* window, select the entries that you want to analyze using the check-boxes next to the entries or with the **Ctrl**- or **Shift**-keys. In this example, select the entry with key "SRR6128338".

2. In the *Experiment presence* panel click on the green dot corresponding to the denovo experiment for the selected entry to open the *Sequence editor* window (see Figure 3).

The genome sequence of the entry with key "SRR6128338" is shown in the *Sequence Editor* panel. The annotation obtained from the wgMLST analysis is presented graphically in the *Sequence Viewer* panel and listed in the *Annotation* panel.

3. Highlight a feature in the feature list of the *Annotation* panel to make the annotation information available in the panel on the right (see Figure 3).
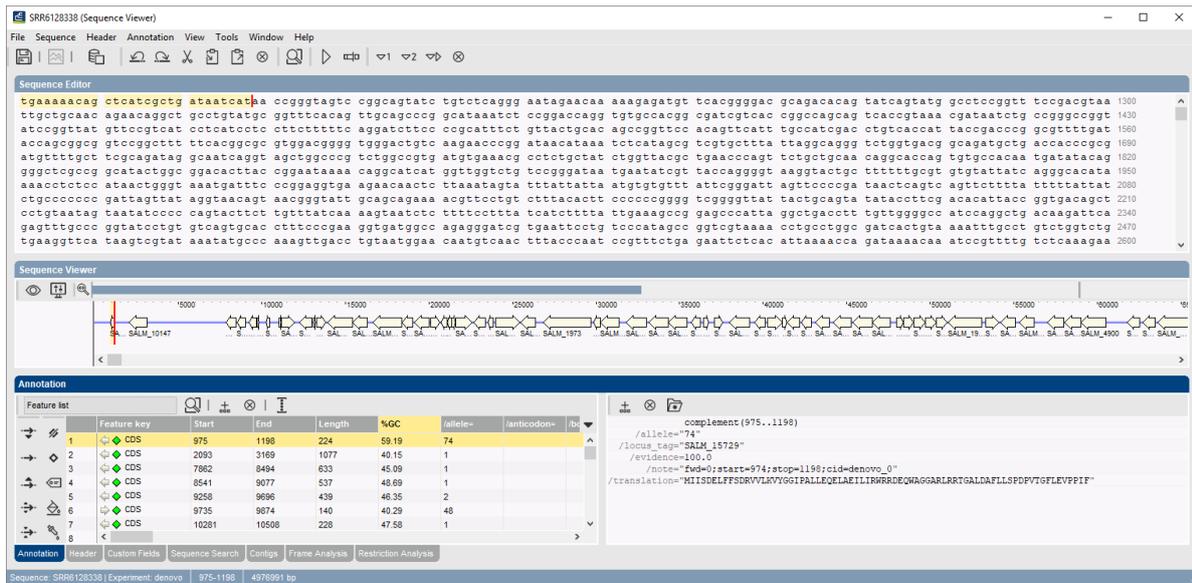
**Figure 3:** The *Sequence editor* window for the selected entry before annotation by Prokka.

4. Close the *Sequence editor* window and select **WGS tools** > **Submit jobs...** ( ▷ ) to call the *Submit jobs* dialog box (see Figure 4).
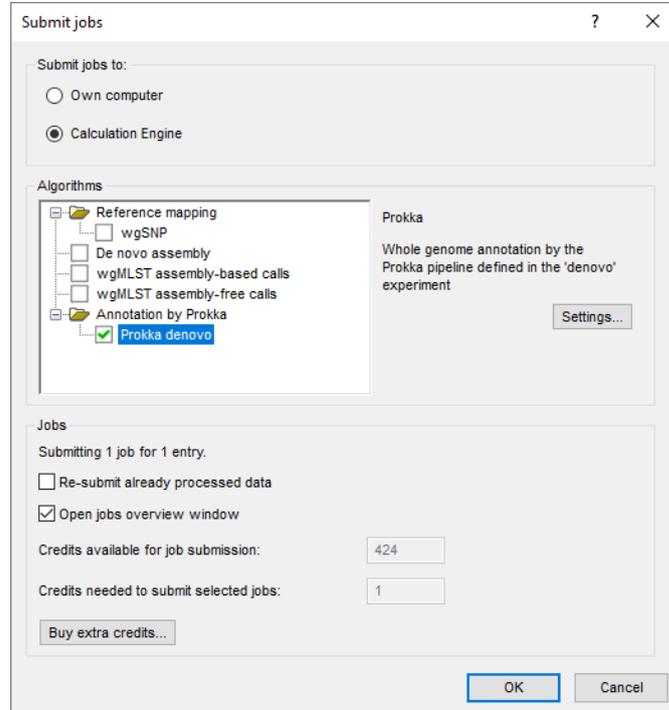


**Figure 4:** The *Submit jobs* dialog box.

5. To run an **Annotation by Prokka** pipeline job on the cloud calculation engine, check the **Calculation Engine** option in the **Submit jobs to** panel and the **Prokka denovo** option in the **Algorithms** panel.

✎ Note that if the selected entry contains sequence data in additional non-reference mapped sequence experiment types in the database these experiment types would also be available under the **Annotation by Prokka** algorithm.

The settings for the Prokka job can be defined by highlighting the job type and pressing <***Settings...***>.

   6. Highlight the **Prokka denovo** job type and press <***Settings...***>.

This action displays the *Prokka settings* dialog box (see Figure 5).



**Figure 5:** The *Prokka settings* dialog box.

Checking ***Force GenBank/ENA/DDJB compliance*** will make the annotations compliant with submission criteria from the GenBank, ENA and DDJB online repositories: add 'gene' features for each 'CDS' feature and a minimum contig length of 200 bp.

By default, Prokka tries to clean up the '/product' names to ensure they are compliant with GenBank/ENA conventions. Checking ***Do not clean up '/product' qualifier annotation*** will prevent this behavior.

The minimum size of a contig to be considered for annotation (***Min. contig length***) can be entered in bp.

In case the corresponding information is already present in the BIONUMERICS database, optionally a ***Genus field***, ***Species field***, ***Strain field*** and/or ***Plasmid field*** can be specified. Entry information contained in these fields will then be included in the annotation, which facilitates later submission to online repositories.

When altering these settings, one can save the updated settings as defaults to the database with ***Save algorithm settings as default***.

   7. Check the ***Force GenBank/ENA/DDJB compliance*** option and leave the other settings as default.

   8. Press the <***OK***> twice to launch the **Annotation by Prokka** pipeline job on the cloud calculation engine.

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.

# 4 Import and analyze Prokka job results

Once the job has been finished (see Figure 6), the results can be imported in the database by selecting **Jobs** > **Get results** (🔍) from the *Job overview* window.
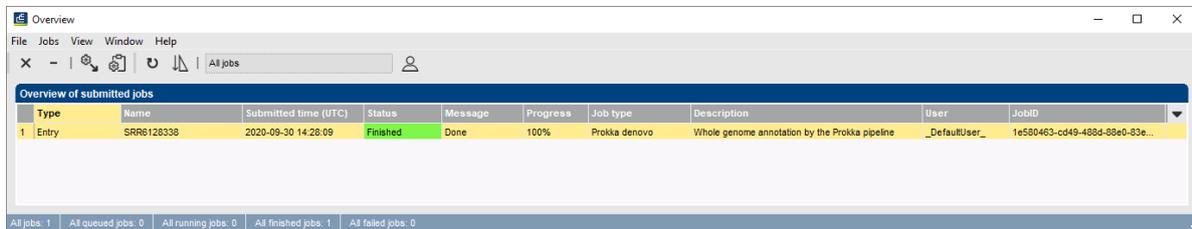


**Figure 6:** The *Job overview* window.

1. When the job is finished, highlight the job and select **Jobs** > **Get results** (🔍) to import the results in the database.

The annotation will be imported in the BIONUMERICS database and saved with the corresponding sequence experiment. A newer Prokka annotation will replace any earlier Prokka annotation on the same sequence, but manually created features or annotation features from other tools will not be overwritten: Prokka features will be added, even if they are defined on exactly the same positions.

2. In the *Experiment presence* panel click on the green dot corresponding to the denovo experiment for the selected entry (i.e. entry with key "SRR6128338") to open the *Sequence editor* window (see Figure 7).



**Figure 7:** The *Sequence editor* window after annotation by Prokka.

The genome sequence of the entry with key "SRR6128338" is shown in the *Sequence Editor* panel. The features detected by Prokka are added to the feature list and can be recognized by the

value "Annotation by Prokka" in the information field /**note=**.

3. Right-click on the **Start** information field and select **Arrange by field** to sort the features based on the start position on the genome. Select a feature in the feature list to make the annotation information available in the panel on the right (see Figure 7).

4. Close the *Sequence editor* window.

The genome sequences can be exported in EMBL or GenBank format.

5. In the *Main* window select **File** > **Export...**.

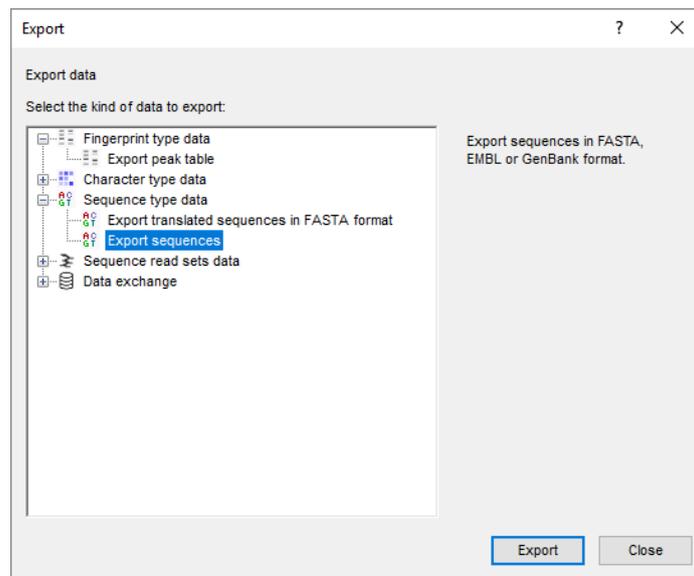The *Export* dialog box opens (see Figure 8).



**Figure 8:** The *Export* dialog box.

6. In the *Export* dialog box select **Export sequences** under **Sequence type data** and <**Export**>.

The *Export sequences* dialog box open (see Figure 9).

7. Browse for the preferred file location, enter a name and press <**Open**>.

8. Select **denovo** as sequence type, **Key** as header field and **GenBank** as output format. Leave the other settings as default and press <**OK**>.

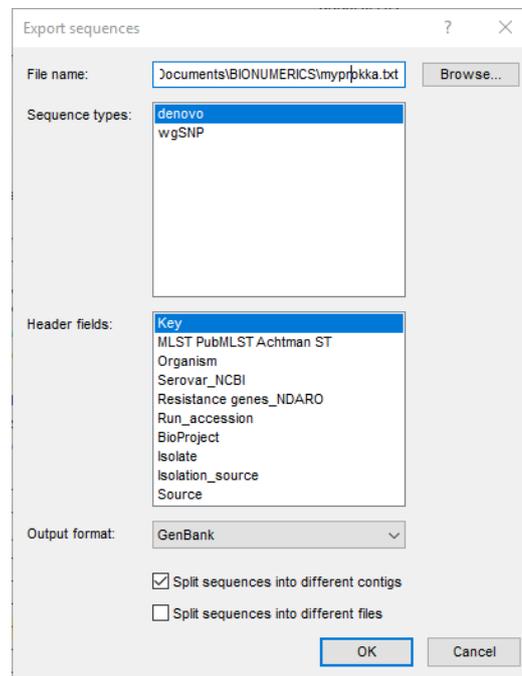The genome sequence with the Prokka annotation will be exported in GenBank format to the provided file location.

**Figure 9:** The *Export sequences* dialog box.

# Bibliography

[1] Torsten Seemann.  Prokka:  rapid prokaryotic genome annotation.  *Bioinformatics*, 30(14):2068–2069, 2014.