BioNumerics Tutorial:

# E. coli functional genotyping: predicting phenotypic traits from whole genome sequences

## 1 Aim

In this tutorial we will screen genome sequences of *Escherichia coli* samples for phenotypic traits using the *E. coli functional genotyping plugin*. This plugin contains public databases for serotype, virulence and resistance prediction, as well as plasmid and prophage detection. An in silico PCR tool is also implemented, making it possible to detect Shiga toxin gene subtypes and virulence genes, mimicking the wet lab PCR.

The different steps are illustrated using the whole genome demonstration database of *Escherichia coli*. This database is available for download on our website (see 2) and contains 60 publicly available sequence read sets of *Escherichia coli* with already calculated de novo assemblies.

## 2 Preparing the database

### 2.1 Introduction to the demonstration database

We provide a **WGS demo database for Escherichia coli** containing sequence read set data links for 60 samples, calculated denovo 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 Escherichia coli** 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 ***Download example databases*** link, located in the lower right corner of the *BioNumerics Startup* window.

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

2. Select the **WGS demo database for Escherichia coli** 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 Escherichia coli** appears in the *BioNumerics Startup* window.

5. Double-click the **WGS demo database for Escherichia coli** in the *BioNumerics Startup* window to open the database.
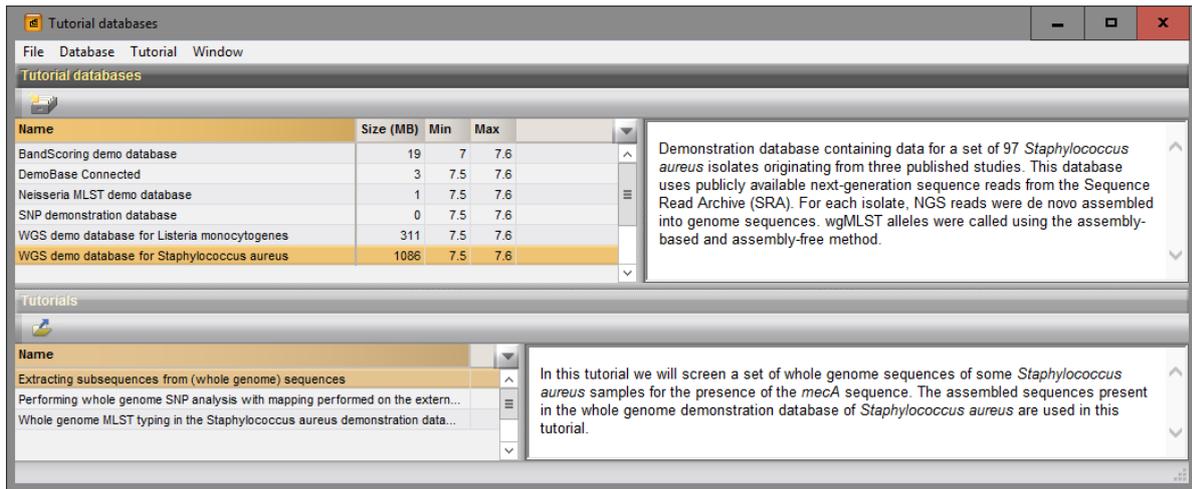
**Figure 1:** The *Tutorial databases* window, used to download the demonstration database.

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

A BioNumerics back-up file of the demo database for Escherichia coli is also available on our website. This backup can be restored to a functional database in BioNumerics.

6. Download the file `WGS_EC.bnbk` file from http://www.applied-maths.com/download/sample-data, under 'WGS demo database for Escherichia coli'.

In contrast to other browsers, some versions of Internet Explorer rename the `WGS_EC.bnbk` database backup file into `WGS_EC.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*** remains selected, an existing database will be overwritten.

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

## 3 About the demonstration database

The **WGS demo database for Escherichia coli** contains data for a set of 60 samples. The sample information, stored in entry info fields (Isolation source, Center Name, etc.) was collected from the publications. Six experiments are present in the demo database and are listed in the *Experiment types* panel (see Figure 2).
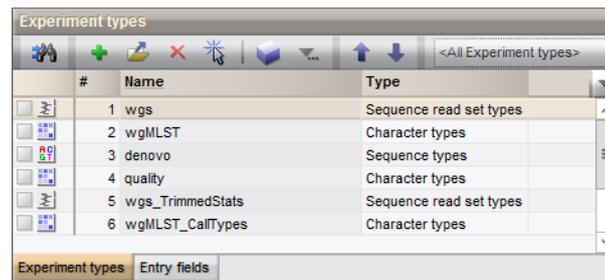
**Figure 2:** The *Experiment types* panel in the *Main* window.

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 default configuration.

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 3).



**Figure 3:** The sequence read set experiment card for an entry.

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** in the default configuration.

The *Sequence editor* window opens, containing the results from the de novo assembly algorithm, i.e. concatenated de novo contig sequences (see Figure 4).
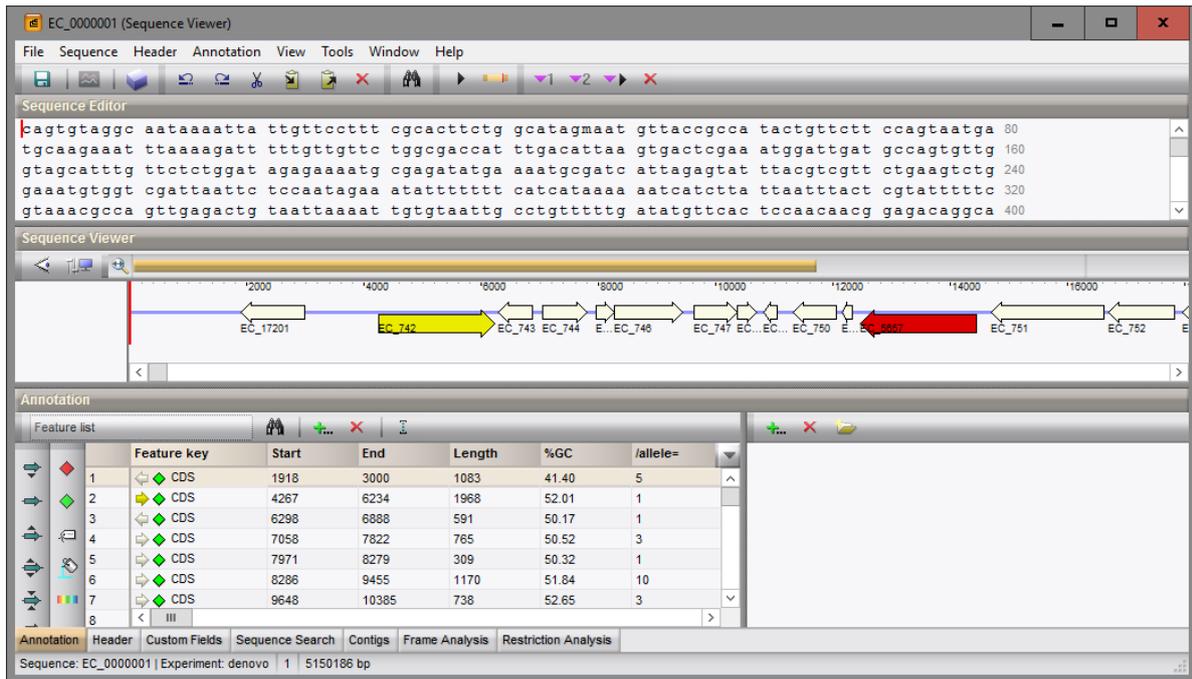
**Figure 4:** The *Sequence editor* window.

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 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.

# 4 Installing the E. coli functional genotyping plugin

1. Call the *Plugins* dialog box from the *Main* window by selecting ***File* > *Install / remove plugins...*** ( ).

The *E. coli functional genotyping plugin* is provided as an *online plugin*. Online plugins are available from the Applied Maths website, from which they can be downloaded and installed in the database in just a few mouse clicks. Since administrator rights are not required for installation of a plugin in the database, online plugins can be easily updated to take advantage of the latest improvements in program code and search data.

2. Select the *Database Functionality tab* in the *Plugins* dialog box and press the <***Add / Update...***> button.

3. Check the box that corresponds to the *E. coli functional genotyping plugin* (see Figure 5) and press <***OK***>.

The actual download of the plugin file can take a couple of minutes, depending on the speed of your internet connection.

The first page of the *E. coli genotyping settings* wizard prompts for some general settings (see Figure 6):
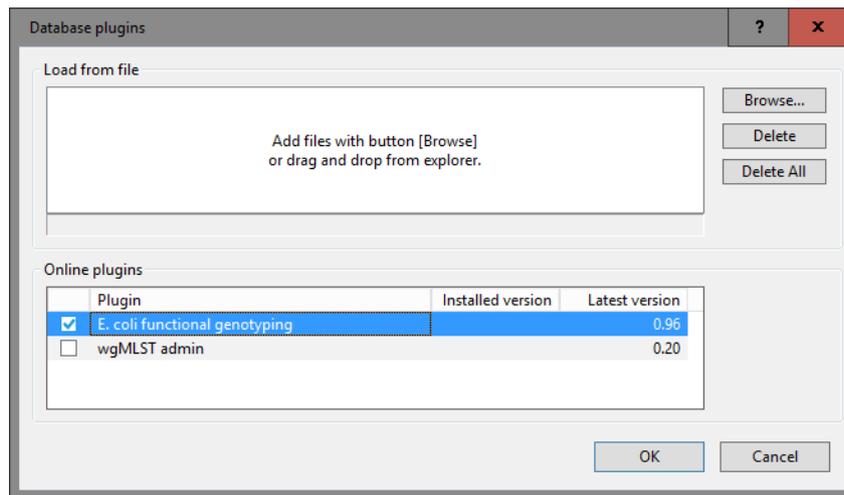
**Figure 5:** Adding an online plugin.

- The ***Sequence experiment type*** that holds the (whole) genome sequences that will be screened.

- The ***Information fields*** that will appear in the report (see 6).

4. In our demonstration database, the assembled sequences are stored in the ***denovo*** sequence experiment. Make sure this experiment is selected from the drop-down list and check the ***Run*** number to include in the report (see Figure 6).
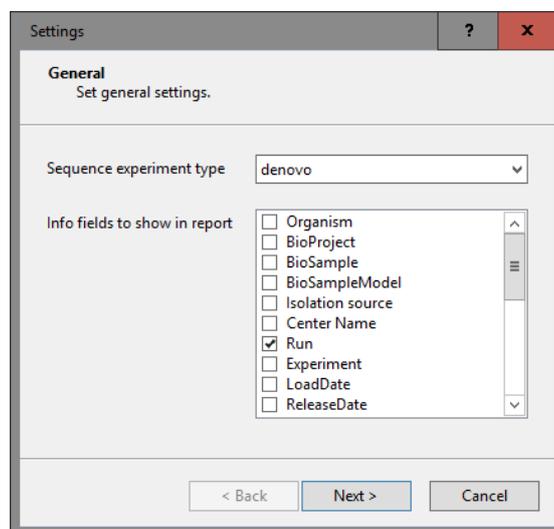


**Figure 6:** The *General settings* wizard page.

5. Press <***Next***>.

The next steps group the settings for each possible search: serotype, pathotype, resistance, virulence, plasmid, prophage, complete plasmid, and in silico PCR search.

The ***BLAST*** search settings become available when checking the ***Determine*** (or ***Find***) check box in each step (see Figure 7 for an example). The ***BLAST settings*** include two thresholds that a BLAST hit should fulfill:

- A ***Minimum sequence identity (%)*** between the subsequence found in the (whole genome) sequence and the sequence in the reference database, expressed as a percentage.

- A ***Minimum length for coverage (%)***, i.e. a minimum overlap between the subsequence found in the (whole genome) sequence and the sequence in the reference database, expressed as a percentage.

With the check box ***Search for gene fragments as well*** enabled in the ***Pathotype***, ***Resistance*** and ***Virulence*** dialogs, sequence fragments (with a maximum of 3 fragments) are also considered.

**Figure 7:** Pathotype settings.

In the *Serotype settings dialog* two additional settings are listed to check the discrimination between the two best hits that passed the BLAST criteria (see Figure 8). Only if the best hit passes these discrimination criteria, the best hit is used for serotype prediction.

- ***Minimum sequence identity for call (%)***: This is the minimum sequence identity the best BLAST hit should have to be considered for serotype prediction.

- ***Minimum discrimination for call (%)***: This is minimum required discrimination between the two best hits in order to predict the serotype based on the best BLAST hit.

**Figure 8:** The *Serotype settings* wizard page.

The discrimination $D$ is calculated as:

$$D = \frac{P_1 - P_2}{100 - P_{min}}$$

with $P_{min}$ = minimum sequence identity percentage specified for BLAST, $P_1$ = sequence identity percentage of best BLAST hit, $P_2$ = sequence identity percentage of second best BLAST hit. When there is only one best hit, $P_2 = P_{min}$.

The last step prompts for the in silico PCR specific settings (see Figure 9):

- *Maximum IUPAC*: Maximum number of allowed IUPAC codes in the subsequence of the (whole genome) sequence. The different possibilities for the ambiguous positions are considered when performing the matching against the sequences in the reference database.

- *Maximum mismatch*: Maximum number of allowed mismatches between the subsequence of the (whole genome) sequence and the sequence in the reference database.



**Figure 9:** The *In silico PCR settings* wizard page.

6. In this tutorial, check all search options and use the default settings in each step.

7. In the last step press *<Finish>* to complete the installation of the plugin.

8. Press *<Exit>* to close the *Plugins* dialog box.

The *E. coli functional genotyping plugin* installs menu items in the main menu of the software under *E. coli* (see Figure 10).

An **In silico PCR overview** character type is added to the *Experiment types* panel. This experiment will summarize the in silico PCR search results. The predicted **In silico PCR** sequences will be stored in the corresponding sequence type experiments (see Figure 11).

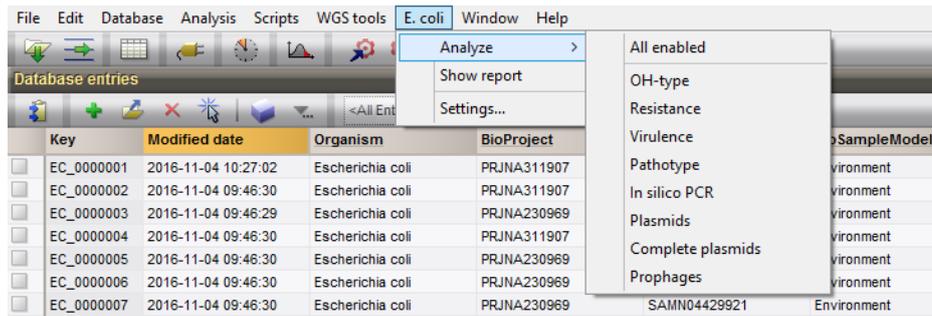The settings specified during installation of the plugin can be called again at any time with *E. coli > Settings...*.

**Figure 10:** New menu-items after installation of the plugin.



**Figure 11:** In silico PCR experiments.

# 5   Screening of entries

The screening can be done on any selection of entries in database.

1. Select a single entry in the *Database entries* panel by holding the **Ctrl**-key and left-clicking on the entry. Alternatively, use the **space bar** to select a highlighted entry or click the ballot box next to the entry.

Selected entries are marked by a checked ballot box (☑) and can be unselected in the same way.

2. In order to select a group of entries, hold the **Shift**-key and click on another entry.

A group of entries can be unselected the same way.

3. Make sure a few entries are selected in the *Database entries* panel of the demonstration database.

Screening for the phenotypic traits can be done for all tools checked in the *Settings* dialog box (***E. coli > Analyze > All enabled***) or for each tool separately (***E. coli > Analyze > ...***).

4. Select ***E. coli  >  Analyze > All enabled*** to screen the selected entries for all enabled traits.

A progress bar appears. The analysis time depends on the number of selected entries. When the analysis is finished, the progress bar disappears. The detected traits for the screened entries are stored in the database.

The **H and O serotypes** that passed the specified call criteria (see Figure 8) and the **pathotypes** that passed the BLAST criteria (see Figure 7) are displayed in the corresponding information fields in the *Database entries* panel (see Figure 12).

| Predicted pathotype | Predicted H serotype | Predicted O serotype |
|---|---|---|
| STEC | H2 | |
| STEC | H2 | |
| STEC | H2 | |
| STEC | H2 | |
| STEC | H2 | |
| STEC | H19 | |
| STEC | H19 | O121 |
| STEC | H19 | O121 |
| STEC | H19 | |
| STEC | H19 | O121 |
| STEC | H19 | O121 |
| STEC | H19 | O121 |
| EPEC | H19 | O121 |
| EPEC | H19 | O121 |
| STEC | H7 | O157 |
| STEC | H7 | O157 |

**Figure 12:** Updated information fields.

The **Resistance**, (**Complete**) **Plasmids**, **Virulence**, **Prophages** character types are created and updated with the predicted traits that passed the BLAST settings.

5. Open a character card (**Resistance**, (**Complete**) **Plasmids**, **Virulence** or **Prophages**) for one of the analyzed entries by clicking on the corresponding green colored dot in the *Experiment presence* panel.

The traits that passed the BLAST criteria are displayed with their identity percentages (see Figure 13 for a few sample character cards).

| EC_0000053 | | | | EC_0000008 | | |
|---|---|---|---|---|---|---|
| **Character** | **Value** | **Mapping** | | **Character** | **Value** | **Mapping** |
| aac(3)-IId | 99.88 | <+> | | FimH | 97.31 | <+> |
| blaTEM-1B | 100.00 | <+> | | iss | 98.37 | <+> |
| mph(A) | 100.00 | <+> | | TraT | 98.63 | <+> |
| | | | | ehxA | 100.00 | <+> |
| | | | | hlyD | 99.92 | <+> |
| | | | | celb | 100.00 | <+> |
| | | | | cif | 99.88 | <+> |
| | | | | eae | 100.00 | <+> |
| | | | | efa1 | 99.71 | <+> |
| | | | | espA | 100.00 | <+> |
| | | | | espB | 100.00 | <+> |
| Press Insert to add character | | | | espI | 98.98 | <+> |

**Figure 13:** Resistance and virulence experiments.

6. Close the character card(s) by clicking in the top left corner of the card.

7. Open the **In silico PCR - overview** character card for one of the analyzed entries by clicking on the corresponding green colored dot in the *Experiment presence* panel.

The **In silico PCR - overview** character card (see Figure 14) lists all in silico PCR sequences that passed the search criteria (see Figure 9).

8. Close the character card by clicking in the top left corner of the card.

The predicted **In silico PCR** sequences are stored in the corresponding sequence type experiments.

9. Clicking on a green colored dot for an in silico experiment opens the *Sequence editor* window displaying the sequence (see Figure 15).

10. Close the *Sequence editor* window.

**Figure 14:** Overview of the in silico PCR search.



**Figure 15:** In silico PCR sequence.

# 6 Reports

    1. Open the genotype report for the selected entries with ***E. coli*** > ***Show report***.

The *Report* window contains a genotype report for each of the selected entries (see Figure 16).

    2. Select another entry in the *Entries* panel to update the results in the *Genotype report* panel.

The data the report was run (***Date***), the Key (***Name***), and information fields checked in the *Settings* dialog box (see Figure 6) are displayed in the *Genotype report* panel, followed by a summary of the results of all analyzed traits.

The discrimination calculated for serotype prediction is displayed next to the predicted serotypes, optionally followed by the second best BLAST hit that passed the BLAST criteria (***alternatives***).

All hits that passed the settings for **Resistance**, **Virulence**, **In silico PCR**, **Pathotype**, (**Complete**) **Plasmids** and **Prophages** screening are listed.
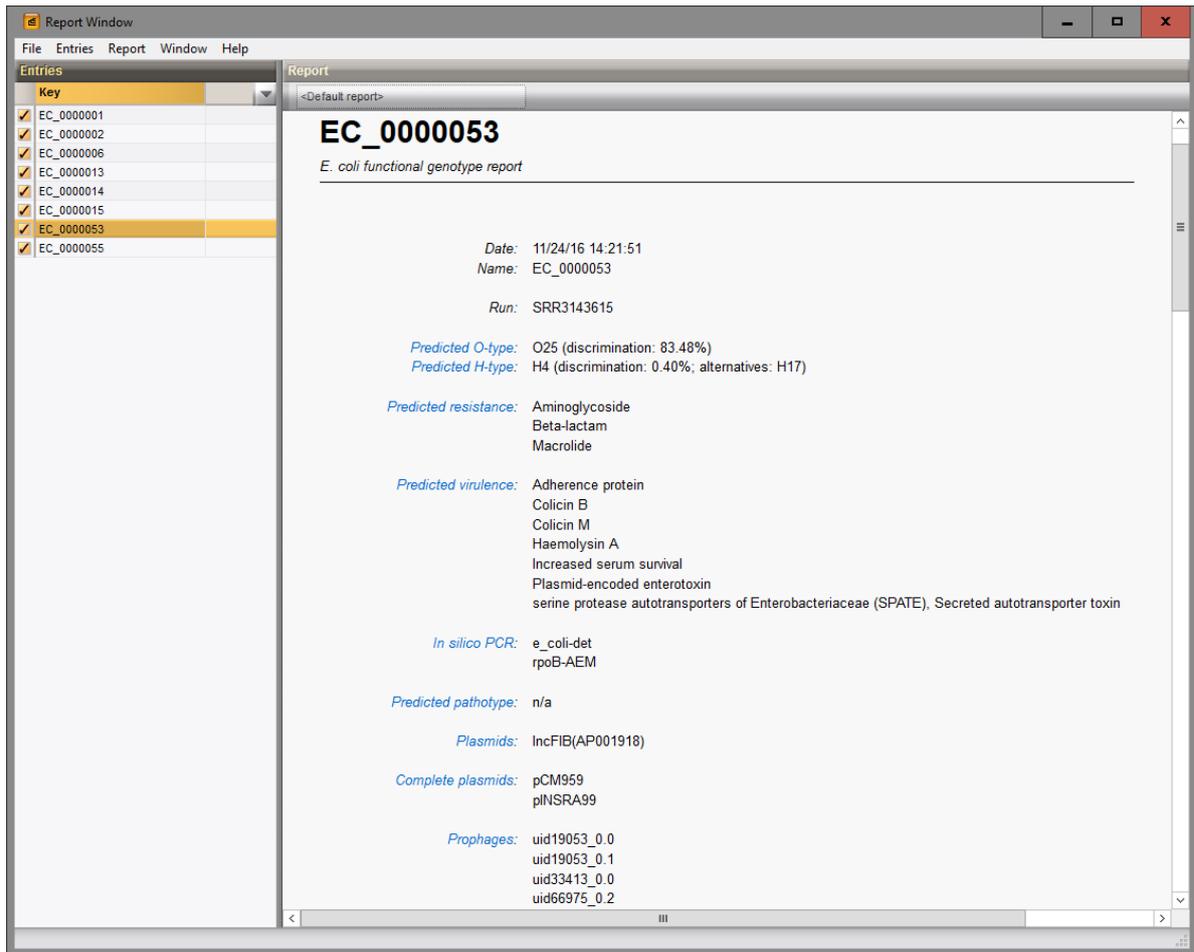
**Figure 16:** Genotype report.

3. Click on a hyperlink of one of the predicted traits to display the detailed BLAST results in the *Genotype report* panel.

Detailed BLAST results include locus identifiers, BLAST similarity scores and descriptive information on the detected genes (see Figure 17). The date the analysis was launched and the version number of the *E. coli functional genotyping plugin* that was used to perform the analysis are indicated.

4. Select ***File*** > ***Exit*** to close the *Report* window.

**Figure 17:** Report details.