

BioNumerics Tutorial:

wgMLST typing in BioNumerics: routine workflow starting from sequence read sets

1 Introduction

This tutorial explains how to prepare your database for wgMLST analysis and how to perform a full wgMLST analysis (de novo assembly, assembly-based and assembly-free calling) in BioNumerics on a routine basis starting from sequence read sets.

The installation of the *WGS tools plugin* requires credentials obtained from Applied Maths. Please make sure you have the credentials ready when following the steps in this tutorial. After installation of the plugin you will notice that importing and analyzing read sets is a very easy and intuitive process.

2 Installation of the plugin

1. Create a new database (see tutorial "Creating a new database") or open an existing database.
2. Call the *Plugins* dialog box from the *Main* window with **File > Install / remove plugins...** .
3. Select the *WGS tools plugin* from the list in the *Applications tab* and press the **<Activate>** button.
4. Confirm the installation of the plugin.

In the first step, the settings for the connection to the **calculation engine** need to be defined. The demanding calculations (i.e. de novo assemblies and allele calling) will be performed on this external calculation engine. The user has the option to use the calculation engine present at the Applied Maths Amazon cloud instance (**Applied Maths cloud**) or to connect to a locally installed instance (**On premises**) (see Figure 1).

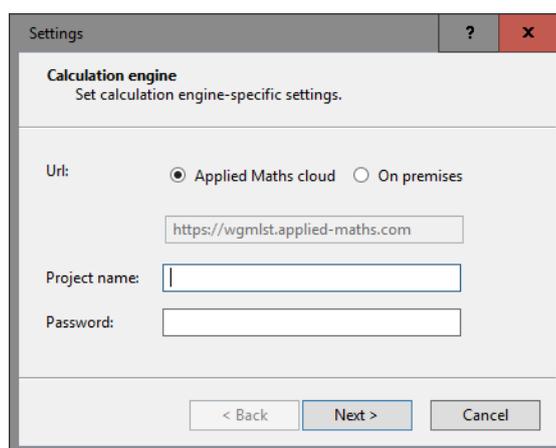


Figure 1: The *Calculation engine* wizard page in the *WGS tools installation* wizard.

5. Select the correct calculation engine resource. In most cases this will be the **Applied Maths cloud** option.
6. Specify the project name as obtained from Applied Maths. The project name is linked with the available credits for a specific account.

7. Specify the password that is used in conjunction with the specified project.
8. Press *<Next>* to proceed to the second step.

The organism (group) can be picked from the drop-down list with available organism schemes. The number of loci is indicated (see Figure 2 for an example).

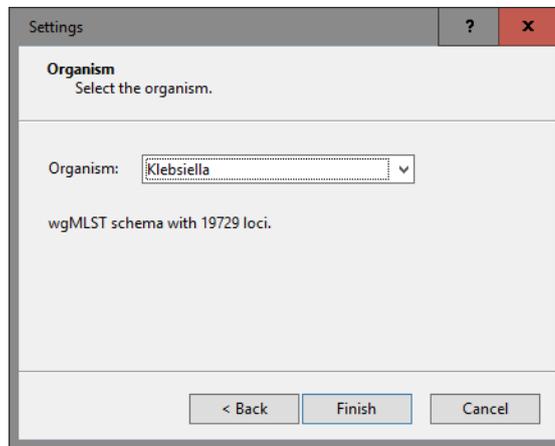


Figure 2: The *Organism* wizard page in the *WGS tools* installation wizard.

9. Press *<Finish>* to start the synchronization with the specified allele database.

The synchronization process can take a couple of minutes, depending on the number of loci and subschemes present in the allele database. A confirmation dialog is displayed when the synchronization has been completed (see Figure 3 for an example).

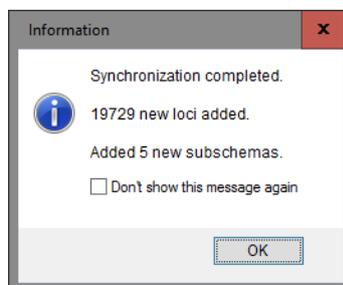


Figure 3: Confirmation of successful installation.

10. Press *<OK>* and *<Exit>* to close the *Plugins* dialog box.



After installation of the plugin, the settings of the *WGS tools* plugin can be accessed with **WGS tools > Settings...**

During installation of the plugin, the **wgMLST** character experiment is created and synchronized with the organism-specific locus scheme. All detected loci and subschemes are added to this experiment.

11. In the *Main* window double-click the character experiment type **wgMLST** in the *Experiment types* panel to call the *Character type* window.
12. Click on the drop-down bar in the toolbar (see Figure 4 for an example).

The views that have been defined at the curator level are synchronized upon installation and are listed. In most databases following views are defined by the curator: the default view **All loci**, the **Core loci** view,

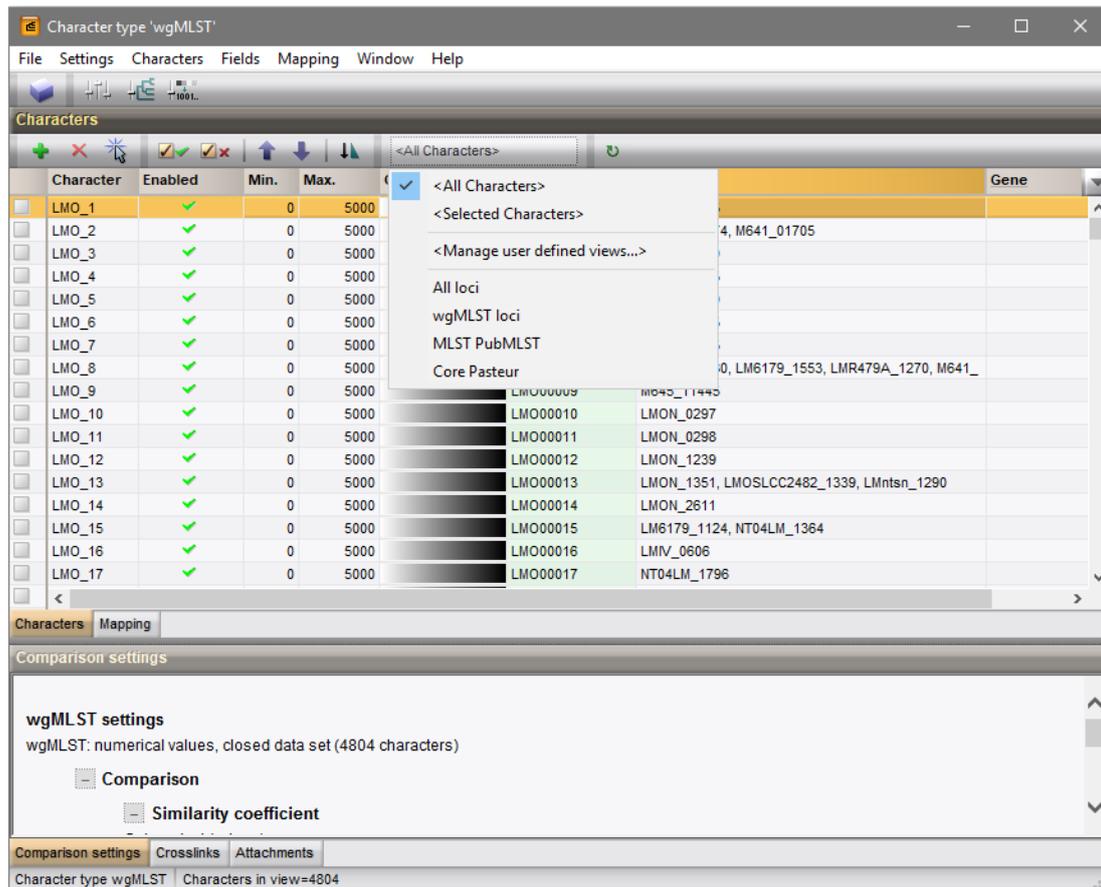


Figure 4: Views defined at the curator side.

the **MLST** view for the traditional seven housekeeping loci, and the **wgMLST loci** view containing all loci except the ones present in the **MLST** view.

13. Select another view from the list to update the set of loci in the *Characters* panel.

The number of loci in the selected view is displayed in the status bar at the bottom of the window.

14. To view all characters again, select **<All loci>** again from the drop-down list.

Besides these curator views, the user can create as many additional local character views as needed and use them as subscheme e.g. for clustering or when inspecting the allele calls for a subset of loci (select *Characters* > *Character Views* > *Manage user defined views*).

15. Close the *Character type* window.

3 Import of sequence read sets

1. Select *File* > *Import...* (📁, **Ctrl+I**) to call the Import tree.
2. Click the +-sign next to the *Sequence read sets data* import option to display the sequence read sets import routines.

Two import routines are listed:

- **Import sequence read sets:** With this option, a multitude of different file types can be imported and

stored inside the database. We do not recommend to use this option since the files might fill up your BioNumerics database quickly and we want to avoid duplication of large data sets.

- **Import sequence read set data as links:** With this option, only the link to the samples is stored in BioNumerics, resulting in a lightweight database. This option is only available after installation of the *WGS tools plugin* and is the preferred option when working with sequence read sets.

3. Make sure the **Import sequence read set data as links** option is selected in the Import tree and press **<Import>** (see Figure 5).

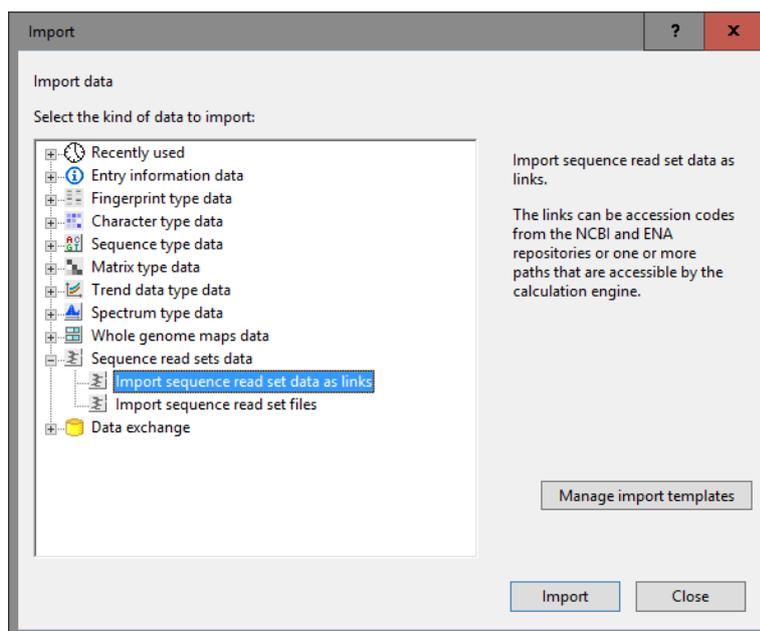


Figure 5: Import sequence read set data as links.

Links to multiple data sources are available, including online and offline data repositories such as: **NCBI (SRA)**, **EMBL-EBI (ENA)**, **Amazon (S3)**, **BaseSpace** or **Local file server** (see Figure 6). Depending on the choice of import, different parameters may be queried in the next steps.

In this tutorial, the import from a local file server is covered. For more information about the other options, please consult the *WGS tools plugin* manual.

4. Select the **Local file server** and press **<Next>**.
5. Press the **<Browse>** button and select your *.fastq or *.fastq.gz files, located on your computer, external drive or on a network location (see Figure 7).

The option **Auto-detect paired-end files** is default checked. This option ensures that the files are checked for the presence of paired-end data. Files that contain paired-end data are recognized by the same file name except for paired-end specific characters: e.g. same name apart from the **_1** or **_2** suffix. Below the file list, a brief summary on the selected files is displayed and updated. This summary indicates how many files of a specific file format were found, and their total file size.

6. Select **<Next>** to go to the next step.

A default import template is listed, parsing the sample names from the file names and linking the sample names to the **Key** field.

7. Press the **<Preview>** button to check the parsing based on the selected template.

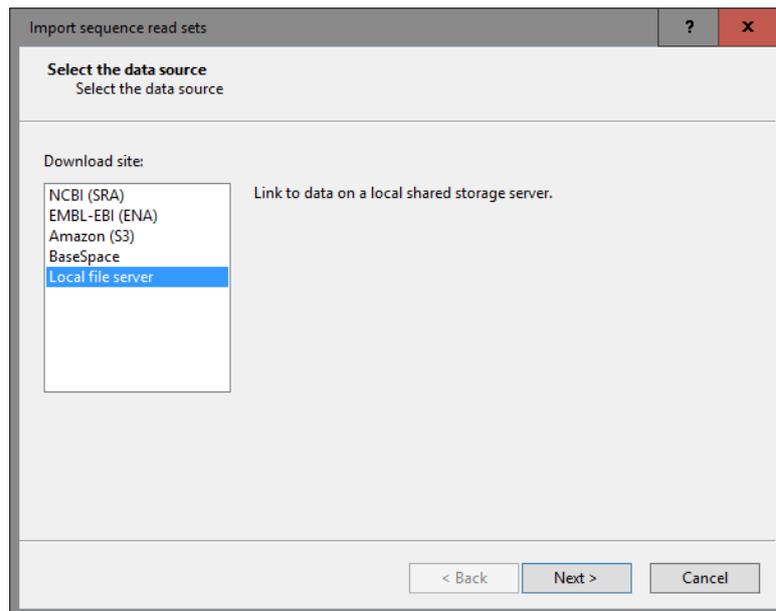


Figure 6: Data sources.

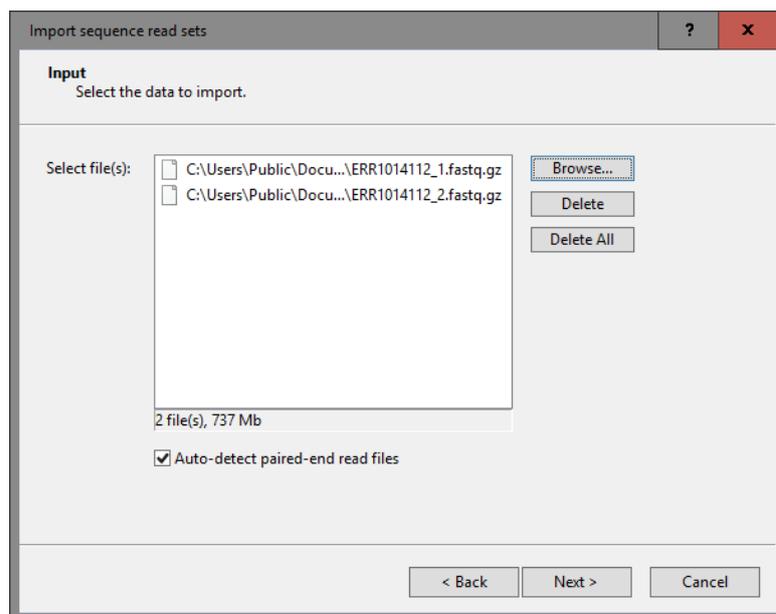


Figure 7: Select sample files.

8. Close the preview.



If the selected import template does not result in a correct parsing of your data, press the **<Edit>** button to change the rules (e.g. link the sample names to an entry information field) or press **<Create new>** to define a new template from scratch.

9. Make sure the **wgs** experiment is selected and click **<Next>** to go to the next step (see Figure 8).

The number of entries that will be created/updated during import is indicated.

10. Click **<Next>**.

In the last step, the wgMLST calculation jobs (de novo assembly, assembly-based and assembly-free calling)

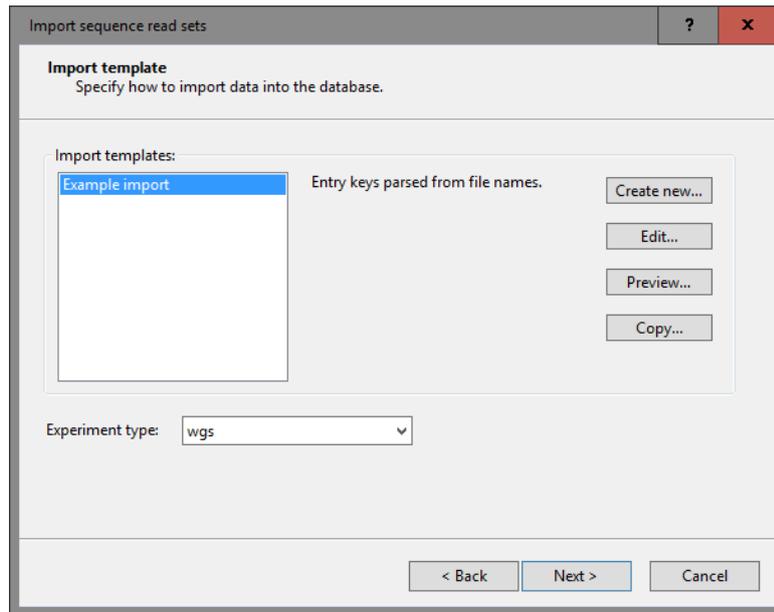


Figure 8: Import template.

can be launched on the imported data links (*Open submit jobs dialog after import*). Note that same dialog can be called from the *Main* window at any time with *WGS tools* > *Submit jobs...* (see 4).

When the *Local file server* option was selected as data source, some basic statistics on the reads can be calculated upon import (*Calculate sequence read set statistics*). Based on the sequence read set statistics bad sequencing runs for which no jobs should be submitted to the calculation engine can be filtered out, saving you credits.

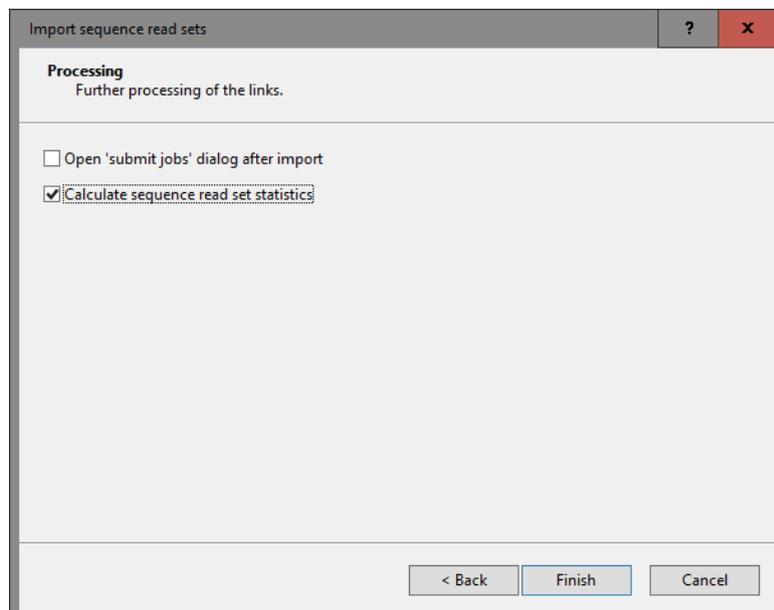


Figure 9: wgMLST calculation jobs.

11. Make sure the *Calculate sequence read set statistics* option is selected, uncheck *Open submit jobs dialog after import* and press <Finish> to start the import of the data links.

Once the import is completed, the entries are created/updated and have one green dot next to it in the column

of the sequence read set experiment type **wgs**.

12. Click on a green colored dot corresponding to the experiment type **wgs**.

The data links are displayed in the *Sequence read set experiment* window. If the option *Calculate sequence read set statistics* was checked in the last step, the statistics are displayed below (see Figure 10). Some statistics are also stored in the **quality** experiment and can be used to filter out bad sequencing results (see 4.1).

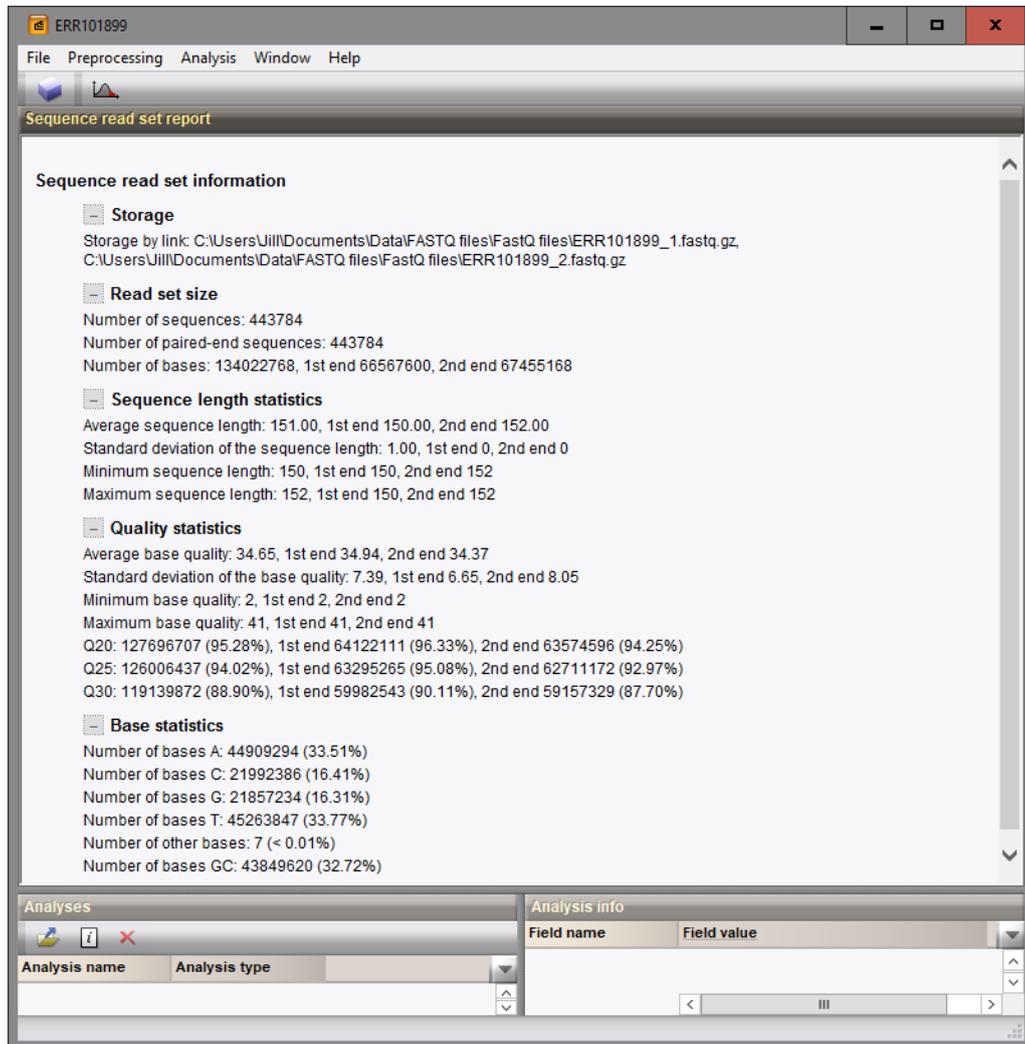


Figure 10: Sequence read set card.

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

4 Submission of jobs

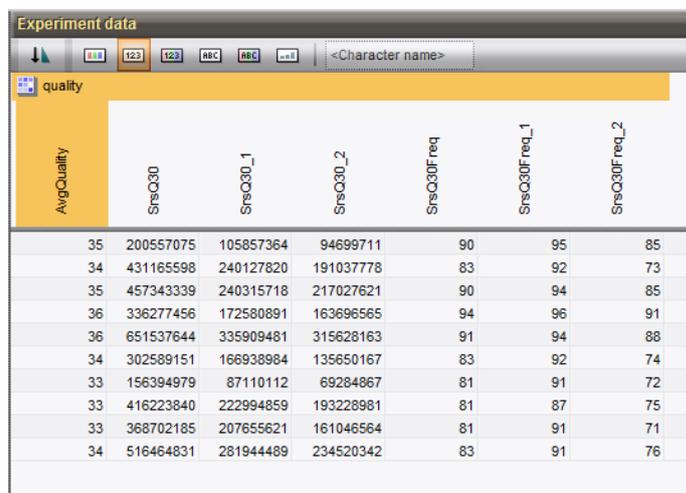
4.1 Check read statistics

Before launching wgMLST jobs on the calculation engine, it is recommended to take a look at the read statistics of the samples to filter out bad sequencing runs, saving you calculation engine credits. These read statistics are calculated when the option *Calculate sequence read set statistics* was checked during import of the read set links (only available when importing reads from a *Local file server*) and are saved in the **wgs**

and **quality** experiments.

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.
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  next to the experiment name **quality** in the *Experiments* panel to display the quality data in the *Experiment data* panel.
4. Select *Characters* > *Show values* () to show the corresponding character values for all entries in the comparison.

The quality values are displayed in the *Experiment data* panel (see Figure 11). The **AvgQuality** is an important indicator for the sequencing quality. The value depends on the sequencing technology used. For Illumina reads for example, the average read quality should be above 30. When samples have Illumina average quality values below 20, these should not be considered for further analysis.



	AvgQuality	SrsQ30	SrsQ30_1	SrsQ30_2	SrsQ30F req	SrsQ30F req_1	SrsQ30F req_2
35	200557075	105857364	94699711	90	95	85	
34	431165598	240127820	191037778	83	92	73	
35	457343339	240315718	217027621	90	94	85	
36	336277456	172580891	163696565	94	96	91	
36	651537644	335909481	315628163	91	94	88	
34	302589151	166938984	135650167	83	92	74	
33	156394979	87110112	69284867	81	91	72	
33	416223840	222994859	193228981	81	87	75	
33	368702185	207655621	161046564	81	91	71	
34	516464831	281944489	234520342	83	91	76	

Figure 11: The quality experiment.

5. Close the *Comparison* window.

To make the distinction between good and bad sequencing run, the quality status (good or bad) can be entered in an entry information field:

6. To create a new entry field, make sure the *Database entries* panel is the active panel in the *Main* window, select *Edit* > *Information fields* > *Add information field...*, specify a name (e.g. **Read quality**) and press **<OK>**. With the option *Edit* > *Information fields* > *Edit field in selection...* (**Ctrl+M**), text can be added to the selected entries (e.g. “Good” or “Bad”).
7. Alternatively, you can opt to permanently remove entries with bad sequencing runs from the database with *Edit* > *Delete selected objects...* ().

4.2 Select jobs

Launching wgMLST jobs on the calculation engine is a very easy process:

8. In the *Main* window, select the entries that you want to analyze using the checkboxes next to the entries or with the **Ctrl-** or **Shift-**keys. Make sure that only samples with good quality reads are included in the selection (see 4.1 to check the good quality samples).

9. Select **WGS tools** > **Submit jobs...** (🔧) to call the *Submit jobs* dialog box.



Alternatively check the **Open submit jobs dialog after import** option in the *Processing* wizard page during import of the data (see Figure 12).

In the *Submit jobs* dialog box you can define which algorithms can be run on the samples (see Figure 12). Three types of jobs are available for **wgMLST**:

- **De novo assembly** to calculate the de novo sequence assembly based on the reads retained after trimming. Two algorithms are available: **Velvet Optimizer** and **SPAdes**. In most cases the **SPAdes** algorithm results in higher quality assemblies.
- **Assembly-based calls** to define the alleles based on a BLAST analysis on the de novo assembled contigs.
- **Assembly-free calls** to define the alleles directly from the reads.



The **Reference mapping** algorithm is used for **wgSNP** analysis and is not covered in this tutorial.

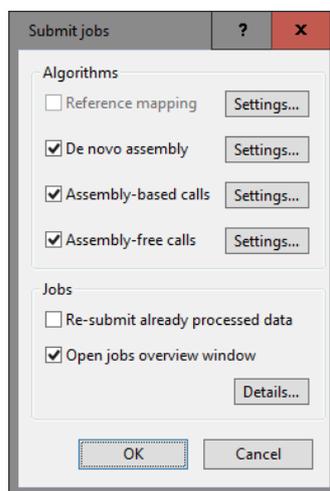


Figure 12: Submit jobs to the calculation engine.

Jobs that already have been submitted and have been imported successfully, will not be relaunched for analysis, unless the check box in front of **Re-submit already processed data** in the **Jobs** part is checked.

Credit costs depend on the job that is submitted: 1 credit for de novo assembly, 3 credits for the Assembly-based calls, and 3 credits for the Assembly-free calls. For more information on the credits press the **<Details>** button.

10. Check the algorithms that you wish to run on the samples, check (and optionally change) the settings, and press **<OK>** to launch the related jobs on the calculation engine.

When links are present to *.fastq or *.fastq.gz files stored on a local hard drive or a local file server a message will pop up asking to upload the files to an Amazon S3 temporary storage (called the **CE Store**), which the calculation engine can access (see Figure 13). Press **<OK>** to start the **CE Store Uploader** (see Figure 14).

4.3 Overview of the jobs

11. By default, the *Calculation engine overview* window will open after submission of the jobs. The same dialog can be called at any time with **WGS tools** > **Jobs overview...** (🔧).

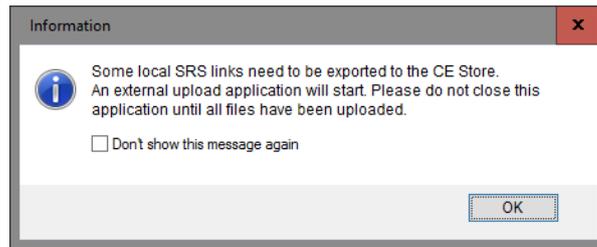


Figure 13: Upload to CE store.

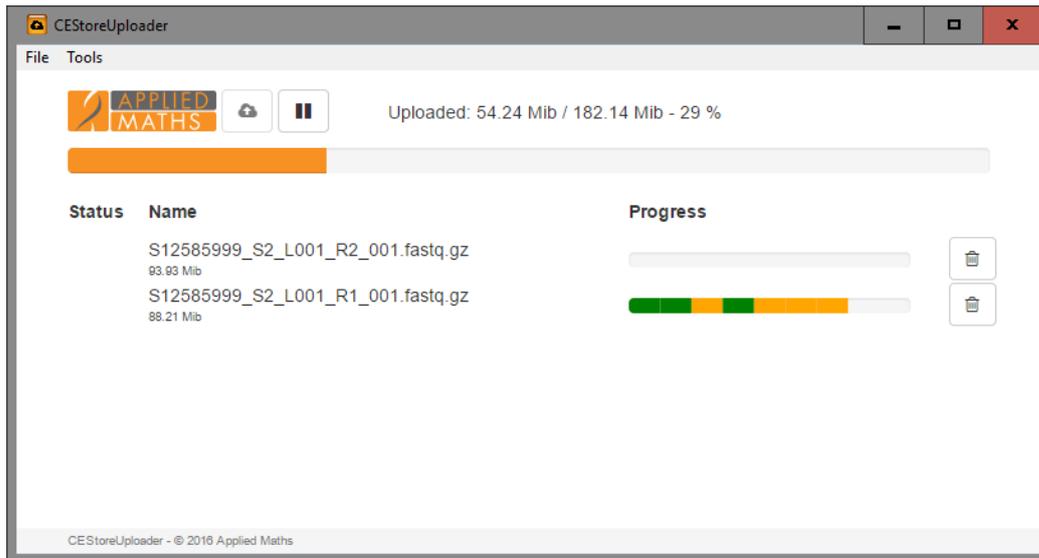


Figure 14: CE Store Uploader.

The *Entry* key, the *Submitted time*, the job *Status*, a *Description* of the job and its *Progress* and much more can be monitored. In the *Message* field, the run comments are displayed in real time (see Figure 15).

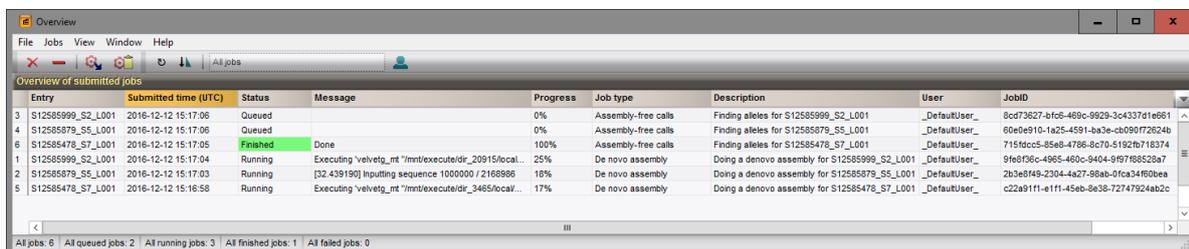


Figure 15: Job overview.

On average, the calculation time for a novo assembly is around **20-30 min**. The Assembly-free calling takes about **5 min**, and the Assembly-based calling is finished after **8 to 9 minutes**.

12. To refresh the overview, press **View > Refresh** (, **F5**).

5 Job results

5.1 Import job results

There are two options available in the *Calculation engine overview* window to import the job results in your BioNumerics database:

1. 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).

The job results can also be imported starting from the entry selection in the *Main* window:

2. Make an entry selection in the *Database entries* panel and select **WGS tools** > **Get results** .

All available job results (for the selected entries) will be imported to the database and linked to their respective entry and experiment type.



The job log files are saved in the *Job log* panel of the *Entry* window. Double-click on an entry in the *Database entries* panel to open the *Entry* window and to consult this information.

Once the results are imported, the corresponding jobs and their underlying data sets are automatically deleted from the calculation engine and as such, from the *Calculation engine overview* window.

Depending on the jobs that were checked in the *Submit jobs* dialog box, following information is stored in the BioNumerics database:

- When the option **De novo assembly** was checked, the results from the de novo assembly algorithm, i.e. concatenated de novo contig sequences with coverage information are stored in the sequence experiment type **denovo**.
- The **wgMLST** experiment contains the allele calls for the detected loci, where the consensus from assembly-based and assembly-free calling - if both jobs were submitted - resulted in a single allele ID. In case multiple allele calls are made and different calls obtained for the same locus, default the lowest common allele ID is retained for these loci (select **WGS tools** > **Settings...** to access this setting in the *wgMLST* tab).
- The character experiment type **quality** contains the quality statistics for the raw data and algorithms that were applied.
- The sequence read set experiment type **wgs_TrimmedStats** contains some data statistics about the reads that were retained after trimming and that were used for the de novo assembly and the assembly-free calling.

5.2 Check job results

The character experiment type **quality** provides insight in the quality of the reads (see 4.1) and the results obtained for the different submitted jobs. The possible presence of low quality reads, assemblies and contaminations can be consulted in a very quick and easy way in the *Comparison* window.

3. 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.
4. Highlight the *Comparisons* panel in the *Main* window and select **Edit** > **Create new object...**  to create a new comparison for the selected entries.
5. Click on the  next to the experiment name **quality** in the *Experiments* panel to display the quality data in the *Experiment data* panel.

6. Select **Characters** > **Show values** (123) to show the corresponding character values for all entries in the comparison.
7. Click on the drop-down list next to the **quality** experiment in the *Experiments* panel to display the default defined character views (see Figure 16).

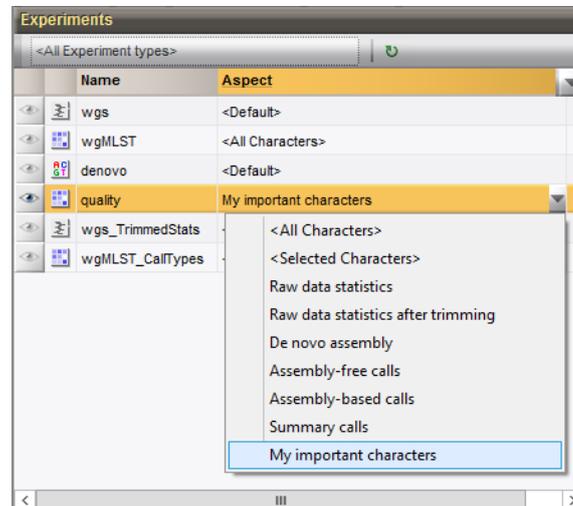


Figure 16: Character views.

The quality parameters are grouped based on the data sets and algorithms and the view can be restricted to each of these groups: raw data statistics (after trimming) (*AvgQuality*, *Srs*, etc.), de novo assembly (*N50*, *NrContigs*, etc.), assembly-free calls (*NrAF*), assembly-based calls (*NrBAF*), and summary calls (*NrConsensus*).

Some parameters are more informative and important than others. A few initial parameters for a first check are listed below:

- **AvgQuality:** the average quality depends on the sequencing technology used. For Illumina reads, the average read quality should be above 30.
- **AvgReadCoverage:** the expected coverage for each base is calculated based on the number of bases in the reads and the expected sequence length. Samples with coverages below 15 should be removed from the analysis. Ideally this number should be above 30.
- **Length:** this length should be close to the length you expect for your organism. Assemblies that are a lot smaller than expected, can be removed from the analysis. For larger lengths, it depends on the cause (contamination or presence of a plasmid).
- **NrAFMultiple:** some loci might have multiple allele hits so a low number is acceptable. If a very high number of multiple allele hits is observed, this indicates the presence of contamination.
- **CorePercent:** this parameter is only present when a core scheme has been defined for the organism. The acceptable range depends on the organism, with typical ranges between 95 and 100. For more diverse organisms, a lower percentage is acceptable, for clonal organisms it is not. This parameter also depends on how strict the core was defined and the diversity of the strains used to define the core. Only very low numbers should be removed without further investigation.

Optionally, you can restrict the view in the *Experiment data* panel to only those parameters that are of

interest to you (see Figure 17 for a custom character quality view):

8. Select the columns in the *Experiment data* panel that you want to include in your custom view while holding the **Ctrl**- key, double-click the **quality** character experiment in the *Experiment types* panel in the *Main* window, select **Characters > Character Views > Manage user defined views...** (**<All Characters>**), press the **<Add>** button, specify a name and select the **Subset based** option.

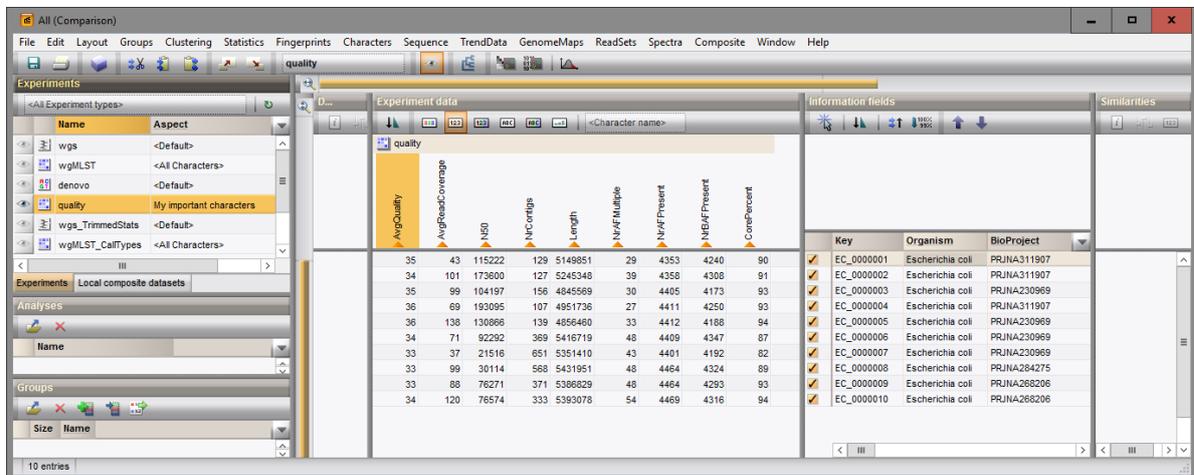


Figure 17: The *Comparison* window displaying a selection of quality parameters.

9. Close the *Comparison* window.

To make the distinction between good and bad job results, this status can be entered in an entry information field:

10. To create a new entry field, make sure the *Database entries* panel is the active panel in the *Main* window, select **Edit > Information fields > Add information field...**, specify a name (e.g. **Job results**) and press **<OK>**. With the option **Edit > Information fields > Edit field in selection...** (**Ctrl+M**), text can be added to the selected entries (e.g. “Good” or “Bad”).
11. Alternatively, you can opt to permanently remove entries with bad job results from the database with **Edit > Delete selected objects...** (**✖**).

6 Follow-up analysis

6.1 Comparison window

A cluster analysis on the **wgMLST** character experiment (or a subscheme thereof) is created in the *Comparison* window or the *Advanced cluster analysis* window.

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.
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 the drop-down list in the **Aspect** column of the **wgMLST** character fields experiment in the *Experiments* panel.

All subschemes defined by the curator in the allele database and the schemes defined by the user (if any) are listed (see Figure 18 for an example). One can very easily switch between the different aspects.

A few analysis tools are highlighted in this tutorial that can be applied on wgMLST data:

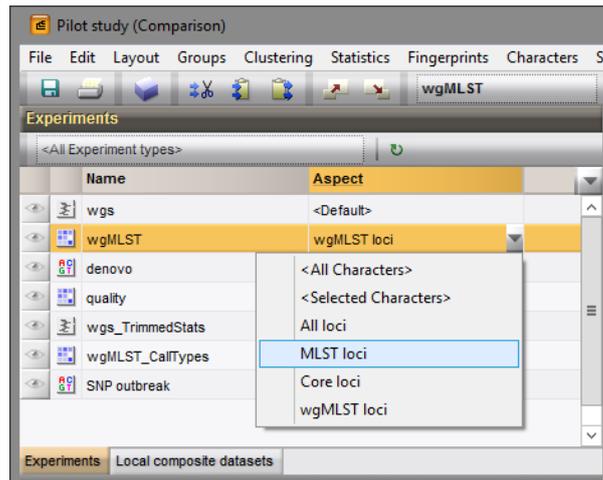


Figure 18: Character views.

6.2 Similarity based clustering

4. Make sure the correct subscheme of the **wgMLST** character experiment that you want to use for your analysis (e.g. **wgMLST loci**, **Core loci**) is selected in the *Experiments* panel.
5. In the *Experiments* panel click on the eye icon (👁) that precedes **wgMLST** to display the values of the selected aspect.
6. In case of closely related isolates select **Clustering** > **Calculate** > **Cluster analysis (similarity matrix)...** and choose the **Categorical (differences)** coefficient from the list (see Figure 19).

The **Categorical (differences)** coefficient treats each different value as a different state, and results in a distance matrix. With the **Scaling factor** one can deal with the hard-coded maximum of 200 that can be calculated for a distance value. Values that make sense are 1, 10 and 100, allowing the correct visualization of maximally 200, 2000 and 20000 different character values, respectively, in a cluster analysis.

7. Press <Next>, choose **Complete Linkage** in the last step and press <Finish>.

When the maximum distance of 200 has been reached, a message is displayed (see Figure 20). To avoid clipping of the dendrogram, repeat the previous steps and increase the **Scaling factor** with 10 or 100.

The resulting dendrogram is displayed in the *Dendrogram* panel and the analysis is stored in the *Analyses* panel. The subscheme that was used is indicated between brackets: e.g. **wgMLST(Core loci)**.

8. The settings used to calculate the dendrogram that is displayed in the *Dendrogram* panel can be called with **Clustering** > **Show information** (📄).
9. To view the number of allele differences on the branches, select **Clustering** > **Dendrogram display settings...** (⚙), and tick the option **Show node information** (see Figure 21).

To trace back the number of different loci from the branches or distance matrix, the displayed values needs to be multiplied with the **Scaling factor** used.

10. The polymorphic loci for the set of samples in the selected scheme can be displayed with **Characters** > **Filter characters** > **Select polymorphic characters...**
11. The information displayed in the *Experiment data* panel can be exported with **Characters** > **Export character table**. The character table will open as a export .csv file in MS Excel.
12. To export the cluster analysis as it appears in the *Comparison* window select **File** > **Print preview...** (🖨), **Ctrl+P**). The *Comparison print preview* window appears.

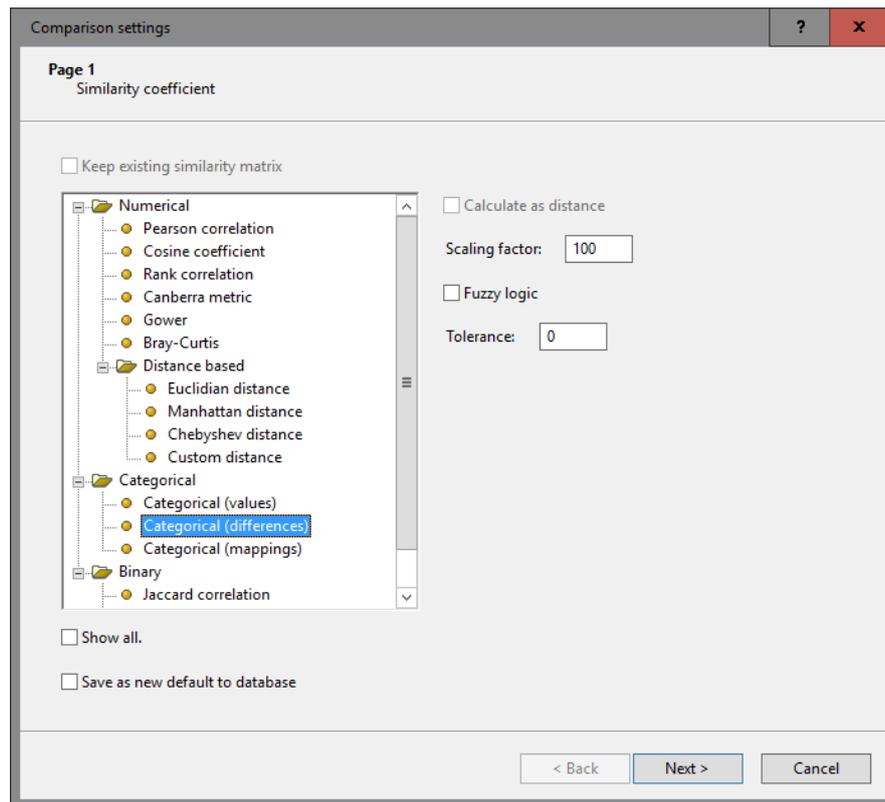


Figure 19: Similarity coefficients.

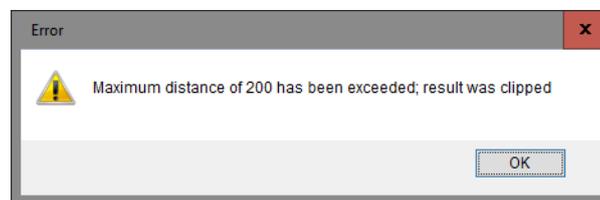


Figure 20: Maximum number.

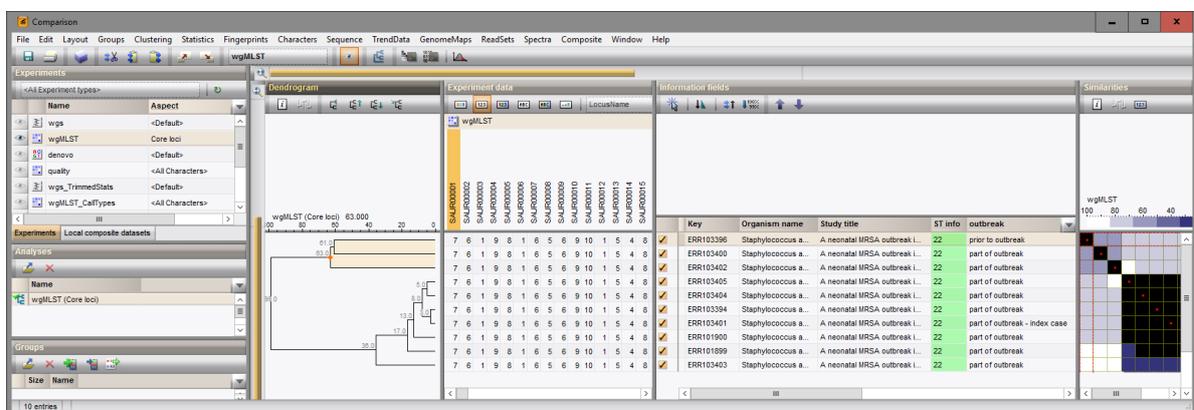


Figure 21: Complete linkage tree based on categorical differences.

More features present in the *Comparison* window are explained in the BioNumerics manual.

6.3 Minimum spanning tree

A minimum spanning tree is calculated in the *Advanced cluster analysis* window which is launched from the *Comparison* window.

13. Select **Clustering** > **Calculate** > **Advanced cluster analysis...** in the *Comparison* window to launch the *Create network wizard*.

The predefined template **MST for categorical data** uses the categorical coefficient for the calculation of the similarity matrix, and will calculate a standard minimum spanning tree.

14. Specify an analysis name, make sure the correct subscheme is selected, select **MST for categorical data**, and press <Next>.



To view and modify the settings of a selected template check the option **Modify template settings for new analysis**.

A MST is now computed in the *Advanced cluster analysis* window (see Figure 22). The *Network panel* displays the minimum spanning tree, the upper right panel (*Entry list*) displays the entries that are present in the tree. The *Cluster analysis method panel* displays the settings used. The analysis is also added to the *Analyses panel* in the *Comparison* window.

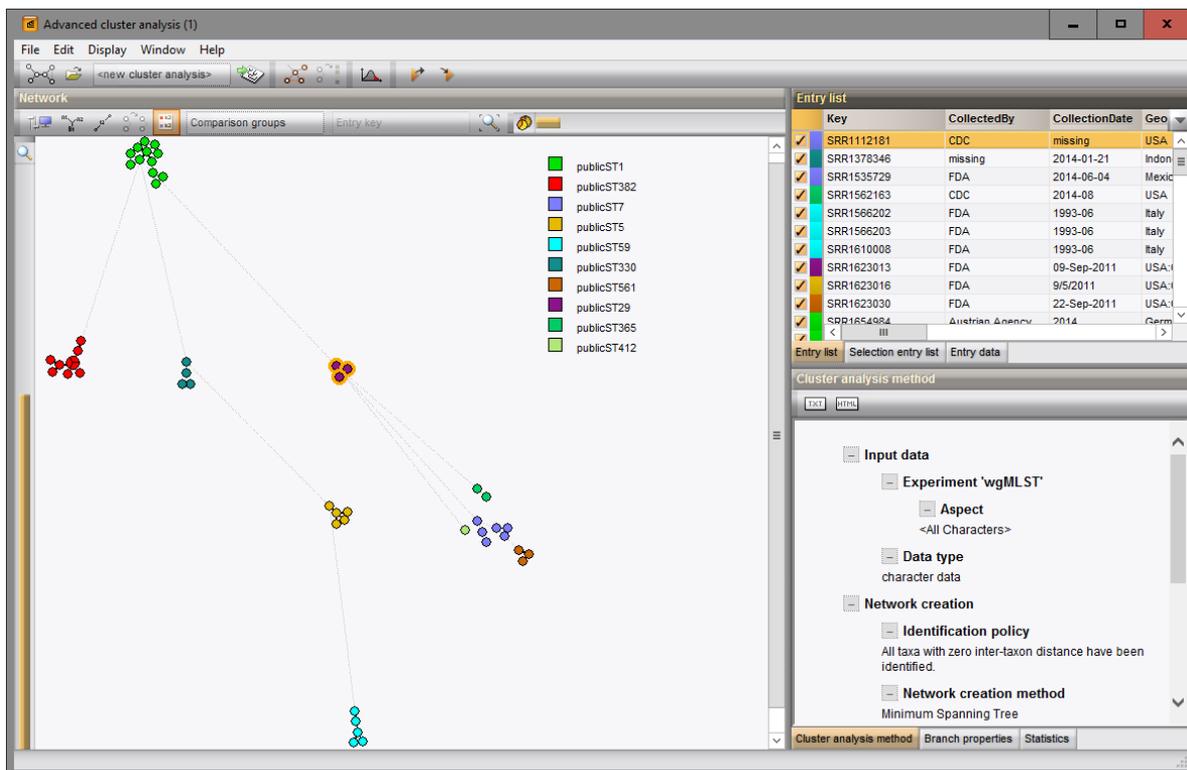


Figure 22: The *Advanced cluster analysis* window.

15. Press  or choose **Display** > **Display settings** to open the *Display settings* dialog box.
16. In the *Branch labels and sizes panel*, you can specify that you want to see the distances between the nodes (i.e. the number of allele differences): check **Show branch labels** and set **Number of digits** to "0".
17. Click <OK> to close the *Display settings* dialog box. The MST is now displayed with branch labels.

More features present in the *Advanced cluster analysis* window are explained in the BioNumerics manual.