



# Tutorial

## MLST-based Outbreak Analysis with Custom Schemes

March 3, 2026

QIAGEN Aarhus A/S · Kalkværksvej 5, 11. · DK - 8000 Aarhus C · Denmark  
[digitalinsights.qiagen.com](https://digitalinsights.qiagen.com) · [ts-bioinformatics@qiagen.com](mailto:ts-bioinformatics@qiagen.com)

Sample to Insight

## MLST-based Outbreak Analysis with Custom Schemes

### Introduction

Multi-Locus Sequence Typing (MLST) is widely used for analysis of microbial genomes, for example in case of an outbreak. There are many curated MLST schemes publicly available from sources such as [PubMLST](#) and [Institut Pasteur](#), but for local outbreaks, it can be beneficial to work with a custom MLST scheme.

This is known as an ad hoc MLST scheme. Using an ad hoc scheme allows for rapid analysis of an ongoing outbreak, by keeping track of allele changes in loci specific to the outbreak.

This tutorial is an introduction to the MLST tools available in *QIAGEN CLC Microbial Genomics Module* by creating an ad hoc MLST scheme and using it for analysis of a local outbreak in a hospital. In some ways, using an ad hoc MLST scheme for outbreak analysis is similar to a SNP-based analysis using variant calling and SNP tree building. For a tutorial including SNP-based analysis, we recommend taking a look at [Typing and Epidemiological Clustering of Common Pathogens](#).

This tutorial covers the following:

- Identifying the best reference(s) for an ad hoc MLST scheme using a publicly available 7-locus MLST scheme.
- Creating an ad hoc MLST scheme.
- Typing reads using an ad hoc MLST scheme.
- Comparing and investigating results from MLST typing.
- Adding typing results to an MLST scheme.
- Creating custom MLST schemes using a workflow.

### Data used in this tutorial

The data used in this tutorial is from [[Hammerum et al., 2015](#)], who investigate a possible hospital outbreak of *Acinetobacter baumannii*. This is an opportunistic pathogen, commonly found in soil and water, but also associated with hospital-acquired infections and is notable for its ability to develop multidrug resistance.

For the sake of time and simplicity, we have selected ten assemblies to use as potential references. These were downloaded using [Download Custom Microbial Reference Database](#). When creating an ad hoc MLST scheme, it is fine to use one or a few closely related references for building the scheme. Otherwise, the references should include as many strains as possible. We advise using 30-50 high quality reference assemblies.


Note that [Create MLST Scheme](#) requires that at least one of the input references has CDS annotations. The data for this tutorial is already annotated. If you wish to use unannotated data as the basis of an MLST scheme, such as a *de novo* assembly, you should first annotate it (e.g., using one of the [Functional Analysis](#) tools available in *CLC Microbial Genomics Module*).

## Prerequisites

For this tutorial, you must be working with *CLC Genomics Workbench 26* and *CLC Microbial Genomics Module 26* or higher. Note that higher versions may produce slightly different results than those shown here. Additionally, external resources like the publicly available MLST Scheme may have been updated to newer versions.

Installing plugins is described in the [CLC Genomics Workbench manual](#).

## General tips

- Throughout this tutorial, we provide links to relevant manual pages, which we recommend exploring for additional details.
- Tools and workflows can be found in the [Toolbox](#), but it is often easier to launch them using [Quick Launch](#) () , found in the top toolbar (shortcut Ctrl+Shift+T or ⌘ +Shift+T on Mac). Quick Launch displays the full Toolbox path, making it easy to identify the location of the tool or workflow if needed.
- The in-built manual can be accessed by clicking the **Help** button on wizards or by selecting the **Help** option under the **Help** menu.
- Within wizards, the **Reset** button can be used to change settings to their default values.
- [Columns in tables](#) can be hidden by unchecking their name in the Side Panel.
- [Columns in tables](#) can be used to sort the rows, by successively clicking on the column name until the desired order (indicated by an arrow next to the column name) is achieved.
- Most of the tools of *CLC Genomics Workbench* require multiple inputs. When many data elements need to be selected, all elements located under a folder can be added by using the options **Add folder contents** or **Add folder contents (recursively)** found in the right-click menu.
- Many data elements produced by *CLC Genomics Workbench* tools have multiple views, indicated as icons in the lower left corner of elements opened in the [View Area](#). Clicking on one of the view icons while pressing the Ctrl (⌘ on Mac) key will open in split view such that both views are visible at the same time. Often, if viewing a table and a graphical representation in split view, selecting entries in the table will highlight them in the graphical representation. The order of the views can be changed using drag and drop, see [Arrange views in View Area](#).

## Import the data

We start by downloading and importing the tutorial data:

1. Download the **tutorial data**.
2. Start *CLC Genomics Workbench*.
3. Import the data using **Standard Import**:
  - (a) Launch **Standard Import** (📁) using **Quick Launch** (🚀).
  - (b) Locate the tutorial data using the **Add files** button and select **Automatic import** (figure 1).

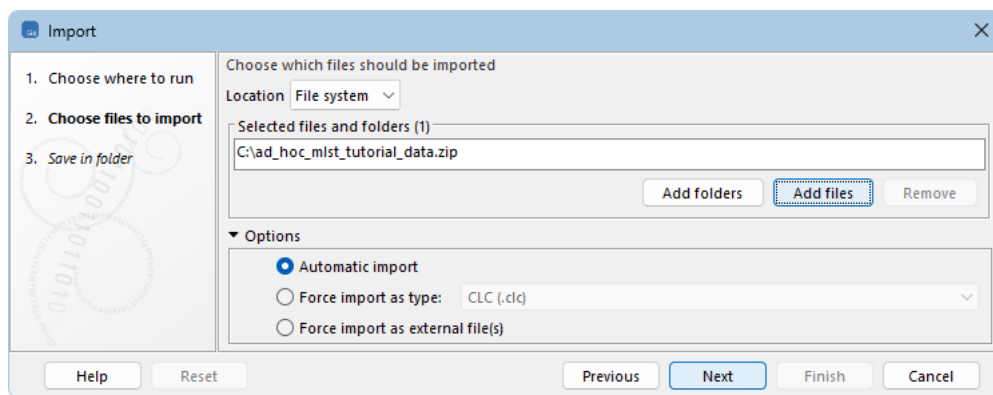


Figure 1: *Standard Import* configured to import the tutorial data.

- (c) In the next step, select a suitable location in the **Navigation Area** to save the imported data and click on **Finish**.

Once the import is completed, an "ad\_hoc\_mlst\_tutorial\_data" folder with two subfolders is visible in the Navigation Area (figure 2).

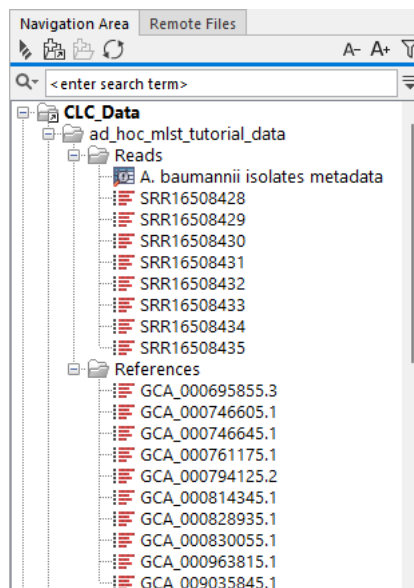


Figure 2: The imported tutorial data in the Navigation Area.

The "Reads" subfolder contains eight sequence lists, one for each set of paired reads, as well as a metadata table. These reads were originally downloaded from SRA and the metadata has been extended with information available in Table 1 of [Hammerum et al., 2015]. The "References" subfolder contains ten reference sequences downloaded from NCBI.

### Determine the best reference(s) for an ad hoc MLST scheme

We will now find the best reference(s) to create a custom MLST scheme for ad hoc analysis of samples.

For this, you must have a user account with **PubMLST**. You must also have registered with the *Acinetobacter baumannii* isolates and *Acinetobacter baumannii* typing databases in your account settings.

If you do not want to create a PubMLST account at this time, you can skip the steps below and jump directly to the section "**Create an ad hoc MLST scheme**", where the best matching reference has already been identified.

### Download a 7-locus MLST scheme

One way to find the best matching reference is by using a 7-locus MLST scheme to type both the samples and the references. First, we will download a publicly available scheme.

1. Launch **Download MLST Scheme** (📄) using **Quick Launch** (🔍).
2. In the first wizard step, "Download settings", select **Acinetobacter baumannii MLST (Pasteur)** under "Scheme to download" and uncheck **Download metadata** to allow for faster download (figure 3).

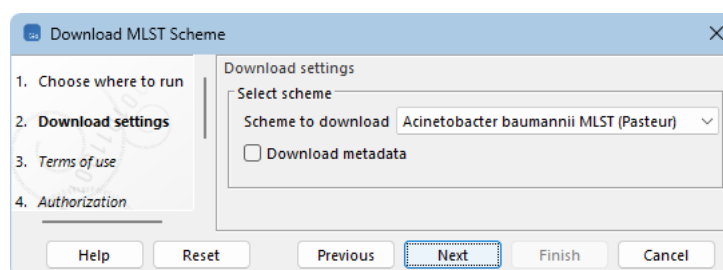


Figure 3: Select the MLST scheme to download.

3. In the next step, "Terms and use", check **I accept these terms**.
4. In the next step, "Authorization", **log in to PubMLST**.
5. In the next step, "Minimum spanning tree parameters", keep the default settings.
6. In the next step, "Result handling", select **Save**.
7. In the last step, "Save location for new elements", click on the **New Folder** (📁) button to create a new subfolder in "ad\_hoc\_mlst\_tutorial\_data" and name it "7-locus MLST scheme".

Choose to save the results in the new subfolder and click on **Finish**.

The download should be quick. We are now ready to type the samples and references using this MLST scheme.

### Type samples and references with the 7-locus scheme

1. Launch **Type with MLST Scheme** (📄) using **Quick Launch** (🔍).
2. In the first wizard step, "Select sequence or sequence list", select the eight samples in the "Reads" subfolder and the ten references in the "References" subfolder (figure 4).

Make sure to check the **Batch** option below the data selection area.

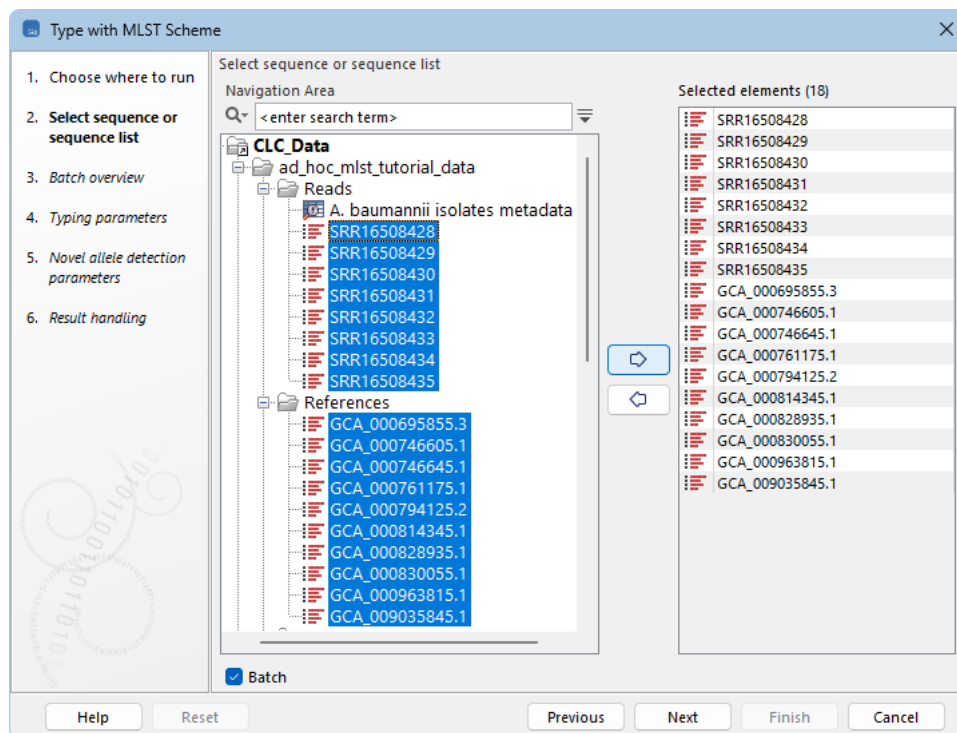


Figure 4: All 18 samples and references are used as input. The batch option must be checked.

3. In the next step, "Batch overview", verify that batch units are as expected. You should see the 18 inputs, which means the tool will run 18 times, once for each input.
4. In the next step, "Typing parameters", select the MLST scheme we just downloaded under "MLST scheme" (figure 5).
5. In the next step, "Novel allele detection parameters", uncheck **Search novel alleles** (figure 6).
6. In the next step, "Result handling", make sure **Create report** is checked.  
Select **Save in specified location** and check **Create subfolders per batch unit**.
7. In the last step, "Save location for new elements", click on the **New Folder** (📁) button to create a new subfolder in "7-locus MLST scheme" and name it "Typing results".  
Choose to save the results in the new subfolder and click on **Finish**.

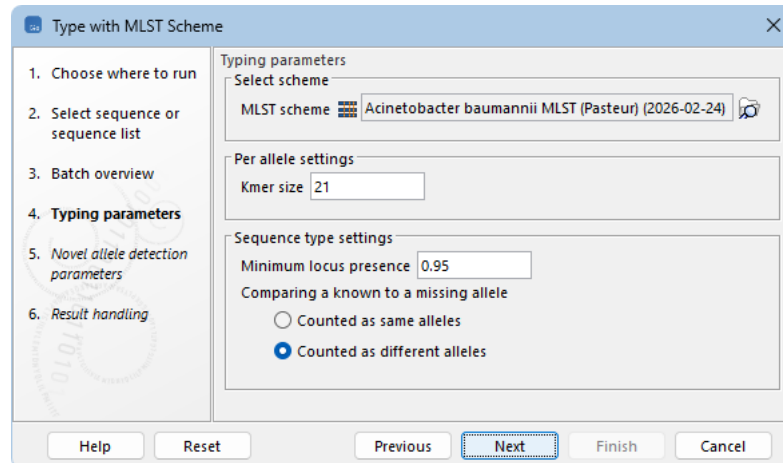


Figure 5: The downloaded *Acinetobacter baumannii* scheme must be used.

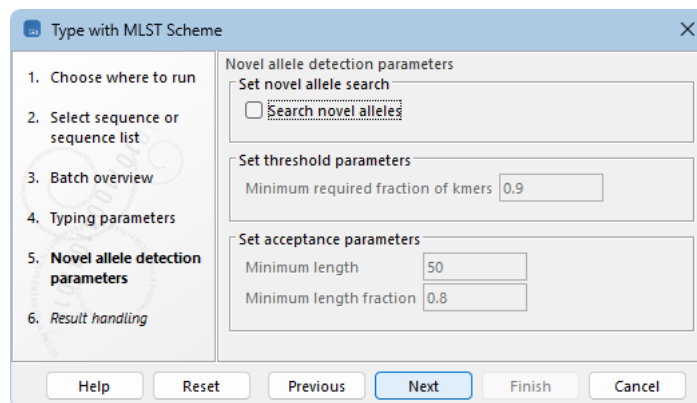


Figure 6: "Search novel alleles" must be unchecked.

The tool will now run 18 times, which may take a while.

Once the tool has finished, we will compare the results. To do so, we will create a combined report:

1. Launch **Combine Reports** (📁) using **Quick Launch** (🚀).
2. In the first wizard step, "Select reports", select the 18 reports as input by right-clicking on the "Typing results" subfolder and selecting **Add folder contents (recursively)** (figure 7).
3. In the next step, "Set contents", keep the default settings.
4. In the next step, "Result handling", select **Save**.
5. In the last step, "Save location for new elements", choose to save the combined report in the "Typing results" subfolder and click on **Finish**.

Open the resulting combined report by double-clicking on it in the Navigation Area. The "Typing result" section shows that of the eight samples, five are typed as ST2, two as ST1 and one as ST158. We can also see that the references consist of various different types (figure 8).

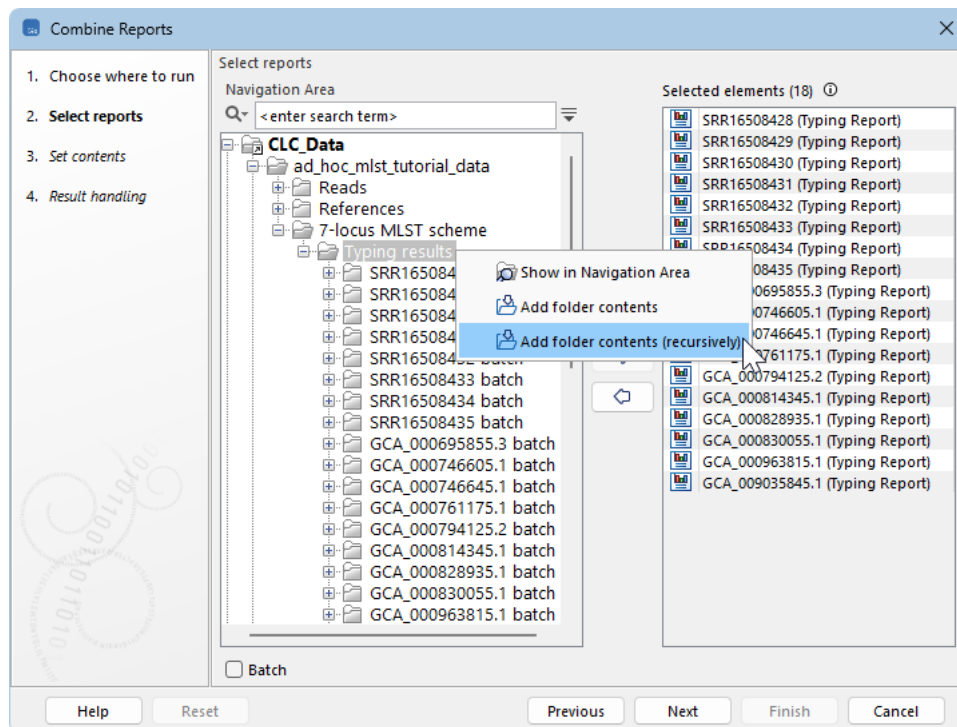


Figure 7: Right-click the top-level folder containing all the outputs from the typing to easily select all reports within that folder.

### 2.1 Typing result

The table is based on 18 samples.

| Sample name                                     | Typing     | Sequence types |
|---|------------|----------------|
| <a href="#">SRR16508428 (Typing Report)</a>     | Conclusive | ST158          |
| <a href="#">SRR16508429 (Typing Report)</a>     | Conclusive | ST2            |
| <a href="#">SRR16508430 (Typing Report)</a>     | Conclusive | ST2            |
| <a href="#">SRR16508431 (Typing Report)</a>     | Conclusive | ST2            |
| <a href="#">SRR16508432 (Typing Report)</a>     | Conclusive | ST1            |
| <a href="#">SRR16508433 (Typing Report)</a>     | Conclusive | ST1            |
| <a href="#">SRR16508434 (Typing Report)</a>     | Conclusive | ST2            |
| <a href="#">SRR16508435 (Typing Report)</a>     | Conclusive | ST2            |
| <a href="#">GCA_000695855.3 (Typing Report)</a> | Conclusive | ST2            |
| <a href="#">GCA_000746605.1 (Typing Report)</a> | Conclusive | ST638          |
| <a href="#">GCA_000746645.1 (Typing Report)</a> | Conclusive | ST79           |
| <a href="#">GCA_000761175.1 (Typing Report)</a> | Conclusive | ST79           |
| <a href="#">GCA_000794125.2 (Typing Report)</a> | Conclusive | ST1            |
| <a href="#">GCA_000814345.1 (Typing Report)</a> | Conclusive | ST464          |
| <a href="#">GCA_000828935.1 (Typing Report)</a> | Conclusive | ST622          |
| <a href="#">GCA_000830055.1 (Typing Report)</a> | Conclusive | ST1            |
| <a href="#">GCA_000963815.1 (Typing Report)</a> | Conclusive | ST1            |
| <a href="#">GCA_009035845.1 (Typing Report)</a> | Conclusive | ST52           |

Figure 8: The "Typing results" section shows the sequence type identified for each sample/reference.

Open the "A. baumannii isolates metadata" table in the "Reads" subfolder. Here, we can see that the column "Pasteur MLST" is in agreement with the results we just obtained, i.e., we get the same sequence types as in [Hammerum et al., 2015].

As our samples contain a mix of sequences types, we will include all references that match these sequence types in our ad hoc scheme, i.e., GCA\_000695855.3 (ST2) and GCA\_000830055.1,



GCA\_000794125.2, and GCA\_000963815.1 (ST1).

None of our included references match ST158, the sequence type of sample SRR16508428. This indicates that for our outbreak analysis, it might be better to obtain and include additional references before building the ad hoc MLST scheme, or to use a more well-established MLST scheme, e.g., one publicly available. For now we will continue our analysis and see how this affects our downstream results.

Alternatively, when no exact match is found to a sequence type, the closest related sequence type can be chosen. We will not show it here, but this can be done by using [Compare MLST Typing Results](#) to evaluate relationships between samples and references.

### Create an ad hoc MLST scheme

Now that we have identified the references most closely related to our samples, we are ready to create an MLST scheme using these references.

1. Launch [Create MLST Scheme](#)  using [Quick Launch](#) .
2. In the first wizard step, "Select contigs or genomic sequences", Select all four ST1 and ST2 genomes to base the MLST scheme on: GCA\_000695855.3, GCA\_000830055.1, GCA\_000794125.2, and GCA\_000963815.1 (figure 9).

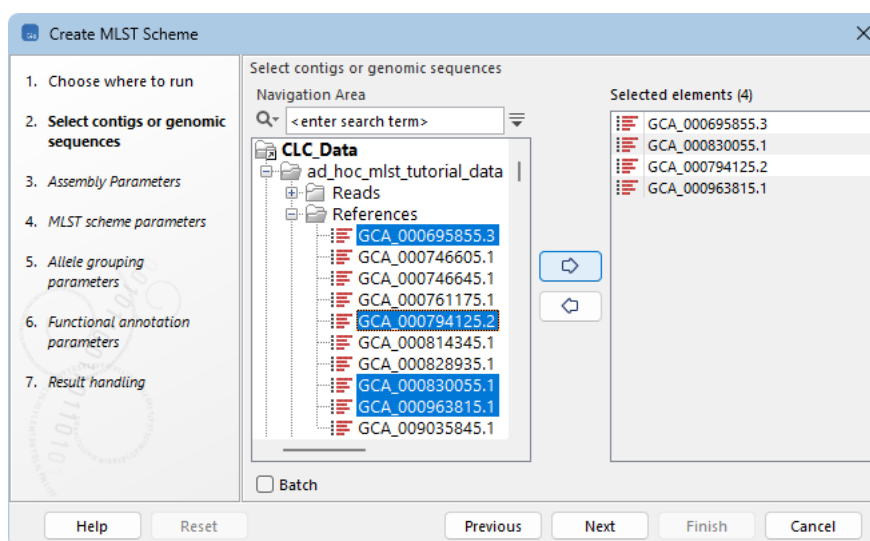


Figure 9: Select the four references typed as ST1 and ST2 as input.

3. In the next step, "Assembly parameters", choose **Each input element is one assembly** under "Assembly grouping" (figure 10).
4. In the next step, "MLST scheme parameters", select the **Whole genome** and **Search alleles before clustering** options to maximize gene detection (figure 11).

Note that if you create schemes from a large number of references, using the Search alleles before clustering option may be quite resource intensive.

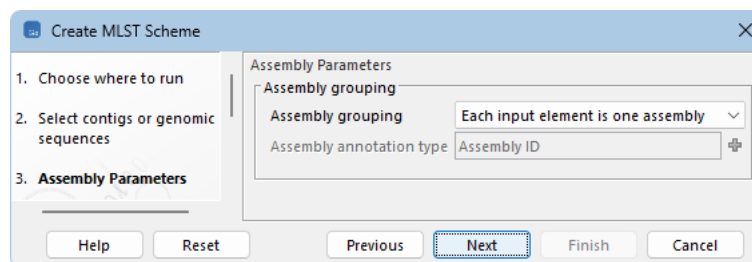


Figure 10: "Each input element is one assembly" must be selected.

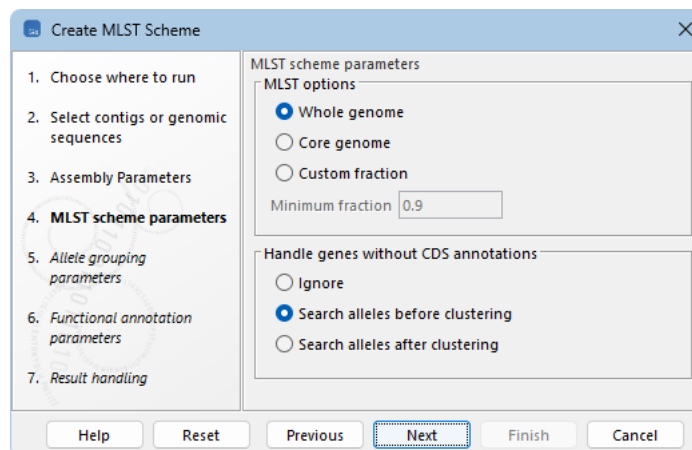


Figure 11: Updated "MLST scheme parameters".

- In the next step, "Allele grouping parameters", keep the default settings.

Note that if you create schemes from a large number of references, it may be necessary to adjust the "Sensitivity" setting under "DIAMOND options" to a less sensitive option to decrease processing time.

- In the next step, "Functional annotation parameters", keep the default settings.

For this tutorial, we will not add any functional annotations, but should you wish to do so, antimicrobial resistance databases and virulence factor databases can be downloaded using [Download Resistance Database](#) (📄).

- In the next step, "Result handling", select **Save**.

- In the last step, "Save location for new elements", click on the **New Folder** (📁) button to create a new subfolder in "ad\_hoc\_mlst\_tutorial\_data" and name it "Create MLST Scheme results".

Choose to save the results in the new subfolder and click on **Finish**.

The tool will now execute. The progress can be monitored under the **Processes** tab in the Toolbox.

Two outputs are created: a report (📄) and an MLST scheme (📄). Feel free to take a look at each of them before moving on. The report contains various information about the input sequences, the number of loci, and the number of alleles, as well as loci/alleles excluded from the scheme and the reason why. The MLST scheme does not yet contain any sequence types, so the only available view is the Allele table view. We will walk through the different views in the scheme after adding sequence types in section "[Add sequence types to the ad hoc MLST scheme](#)".

## Generate and interpret typing results

We will now type our samples using the MLST scheme we just created.

### Type reads using the ad hoc MLST scheme

1. Launch **Type with MLST Scheme** (🚀) using **Quick Launch** (🔍).
2. In the first wizard step, "Select sequence or sequence list", select the eight samples in the "Reads" subfolder.  
Make sure to check the **Batch** option below the data selection area.
3. In the next step, "Batch overview", verify that batch units are as expected. You should see the eight samples, which means the tool will run eight times, once for each input.
4. In the next step, "Typing parameters", select the ad hoc scheme we just created from the "Create MLST Scheme results" subfolder under "MLST scheme" and set **Minimum locus presence** to 0.75 (figure 12).

The **Minimum locus presence** option determines the threshold for whether a sample can be successfully typed. For cgMLST typing, 0.95 is suitable. When working with a wgMLST scheme, this threshold should however be lowered.

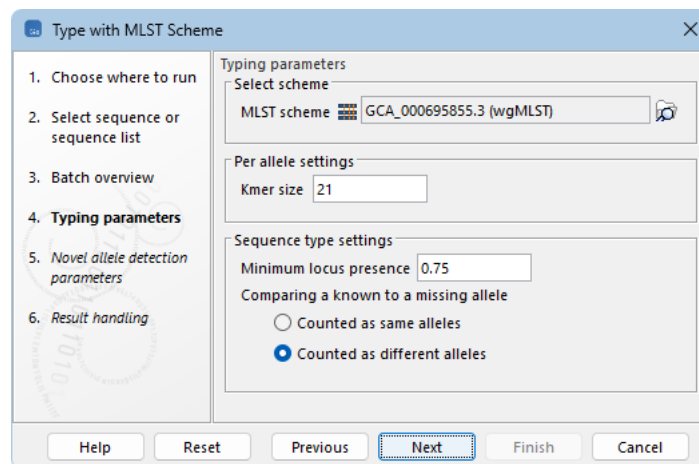


Figure 12: The ad hoc MLST scheme must be used and the "Minimum locus presence" must be lowered.

5. In the next step, "Novel allele detection parameters", check **Search novel alleles** (figure 13).  
Searching for novel alleles allows the tool to detect allelic variants not yet present in the scheme, which is useful when analyzing new or evolving strains. However, this can make results less stable: samples may frequently be flagged as "novel", and results can change later if those alleles are added to the scheme. For well-established MLST schemes where long-term consistency is important, it is often better to disable novel allele detection to keep typing results more stable.
6. In the next step, "Result handling", make sure **Create allele sequence list** and **Create report** are checked.  
Select **Save in specified location** and check **Create subfolders per batch unit**.

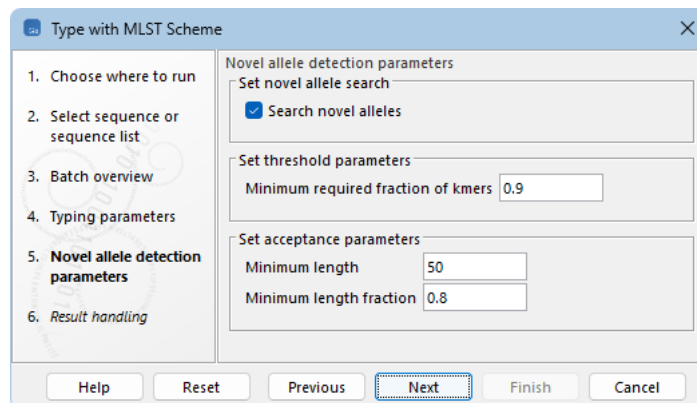


Figure 13: "Search novel alleles" must be checked.

7. In the last step, "Save location for new elements", click on the **New Folder** (📁) button to create a new subfolder in "ad\_hoc\_mlst\_tutorial\_data" and name it "Ad hoc typing results". Choose to save the results in the new subfolder and click on **Finish**.

The tool outputs a **Typing Result** (📄), a **Typing Report** (📄), and an **Allele sequence list** (☰) for each sample. We will explore these outputs below, first comparing the successfully typed samples, and then taking a closer look at the failed sample.

### Compare typing results in a combined report

First, we will investigate the typing reports. We will combine the reports in the same way as we did in section "[Type samples and references with the 7-locus scheme](#)".

1. Launch **Combine Reports** (📄) using **Quick Launch** (🚀).
2. In the first wizard step, "Select reports", select the eight reports as input by right-clicking on the "Ad hoc typing results" subfolder and selecting **Add folder contents (recursively)** (figure 14).

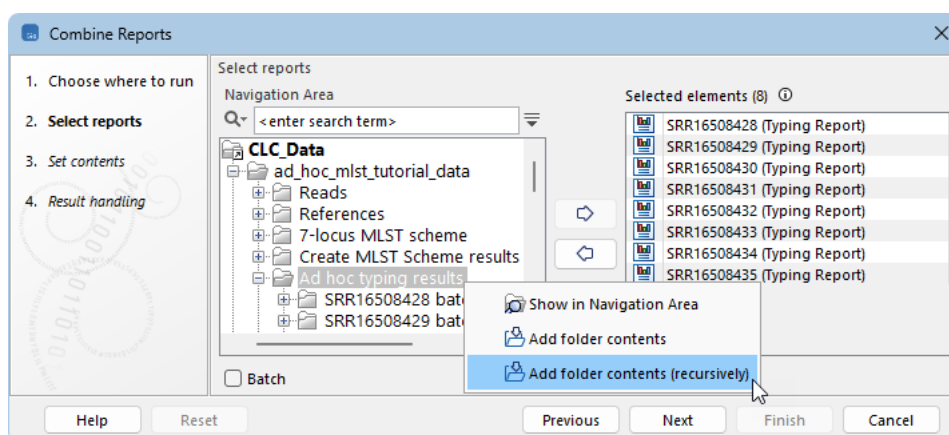


Figure 14: Right-click on the top-level folder containing all the typing results to easily add all reports within that folder.

3. In the next step, "Set contents", keep the default settings.

4. In the next step, "Result handling", select **Save**.
5. In the last step, "Save location for new elements", choose to save the combined report in the "Ad hoc typing results" subfolder and click on **Finish**.

Open the resulting combined report by double-clicking on it in the Navigation Area. From the "Typing result" section, we can see that seven of the samples are successfully typed as "Novel". Samples will always be typed as "Novel" when there are no known sequence types in the MLST scheme (figure 15).

We also see that for one of the samples (SRR16508428), the typing is considered "Not possible". This is the same sample for which we did not have a matching ST158 reference genome in section "[Determine the best reference\(s\) for an ad hoc MLST scheme](#)".

From the "Typing summary" section, we see that this is because the locus presence falls below the threshold of 0.75 that we set in step 4 (figure 15).

### 2.1 Typing result

The table is based on 8 samples.

| Sample name                 | Typing       | Sequence types |
|-----------------------------|--------------|----------------|
| SRR16508428 (Typing Report) | Not possible | -              |
| SRR16508429 (Typing Report) | Conclusive   | Novel          |
| SRR16508430 (Typing Report) | Conclusive   | Novel          |
| SRR16508431 (Typing Report) | Conclusive   | Novel          |
| SRR16508432 (Typing Report) | Conclusive   | Novel          |
| SRR16508433 (Typing Report) | Conclusive   | Novel          |
| SRR16508434 (Typing Report) | Conclusive   | Novel          |
| SRR16508435 (Typing Report) | Conclusive   | Novel          |

### 2.2 Typing summary

The table is based on 8 samples.

| Sample name                 | Loci identified | Novel alleles identified | Loci not identified | Locus presence (%) | Estimated sample coverage |
|-----------------------------|-----------------|--------------------------|---------------------|--------------------|---------------------------|
| SRR16508428 (Typing Report) | 544             | 317                      | 1,706               | 24.18              | 42.34                     |
| SRR16508429 (Typing Report) | 1,716           | 90                       | 534                 | 76.27              | 32.89                     |
| SRR16508430 (Typing Report) | 1,817           | 49                       | 433                 | 80.76              | 20.30                     |
| SRR16508431 (Typing Report) | 1,821           | 48                       | 429                 | 80.93              | 41.29                     |
| SRR16508432 (Typing Report) | 1,707           | 72                       | 543                 | 75.87              | 30.91                     |
| SRR16508433 (Typing Report) | 1,695           | 74                       | 555                 | 75.33              | 25.86                     |
| SRR16508434 (Typing Report) | 1,814           | 28                       | 436                 | 80.62              | 44.77                     |
| SRR16508435 (Typing Report) | 1,808           | 29                       | 442                 | 80.36              | 28.43                     |
| Minimum                     | 544.00          | 28.00                    | 429.00              | 24.18              | 20.30                     |
| Median                      | 1,762.00        | 60.50                    | 488.00              | 78.31              | 31.90                     |
| Maximum                     | 1,821.00        | 317.00                   | 1,706.00            | 80.93              | 44.77                     |
| Mean                        | 1,615.25        | 88.38                    | 634.75              | 71.79              | 33.35                     |
| Standard deviation          | 436.25          | 94.94                    | 436.25              | 19.39              | 8.71                      |

Figure 15: Sample SRR16508428's locus presence is less than the threshold of 0.75.

Getting a "Not possible" typing result when using a well-established MLST scheme may indicate that the sample is not the same species as the MLST scheme. In this case, however, we know the sample is *Acinetobacter baumannii* and we know the scheme is an ad hoc scheme, made from only two sequence types. Therefore, the "Not possible" result corroborates our previous assumption that for our outbreak analysis, it might be better to obtain and include additional references before building the MLST scheme, or to use a well-established MLST scheme, e.g., one publicly available.

Other ways to handle failed typing results is to lower the minimum locus presence (figure 12) or lower the minimum required fraction of kmers (figure 13) to potentially find more novel alleles.

### Investigate and verify a failed typing result

As mentioned above, sample SRR16508428 could not be successfully typed. To investigate this result, we will now take a closer look at the **outputs from Type with MLST Scheme** in the "SRR16508428 batch" subfolder of "Ad hoc typing results".

- First, open the Typing Report ( (figure 16):

Here, we can see much of the same information as available in the combined report that we created, but we also see a coverage distribution plot. From the plot and the **Estimated sample coverage**, we see that although the frequencies are overall lower than for other samples and there is a relatively high peak at 0, the coverage seems evenly distributed. This corroborates that the low locus presence is indeed due to these loci not being present in the MLST scheme, as opposed to a sequencing depth issue for the whole sample.

#### 2.1 Typing statistics

|                           |       |
|---------------------------|-------|
| Loci identified           | 544   |
| Novel alleles identified  | 317   |
| Loci not identified       | 1,706 |
| Locus presence (%)        | 24.18 |
| Estimated sample coverage | 42.34 |

#### 2.2 Coverage distribution

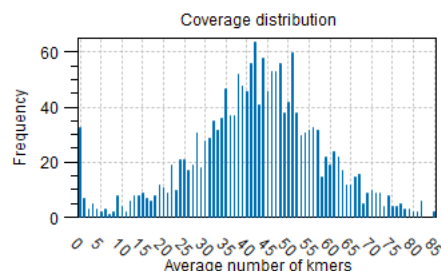

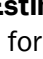



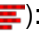
Figure 16: *Typing Report* for sample "SRR16508428".

- Second, open the Typing Result ():

The Allele Table view () shows all the loci for which at least one kmer was found, with the best match allele called or novel allele found for each. Select all columns in the **Side Panel** to the right of the view, to see more information on the kmers found. The "Fraction of kmers" must be 1.0 for an allele to pass.

If we **filter the table** for "Status = Fail", we can see that many of the loci have a relatively high fraction of kmers found as well as total and average kmer counts. This again indicates that we are indeed working with something that is very closely related to our scheme, but that it is just slightly too dissimilar to call novel alleles.

The Novel Allele Table view () contains a sequence list of any novel alleles found during the typing. From here, the sequences can be extracted for further analysis, if desired.



- Third, open the Alleles sequence list ():

This sequence list contains all the alleles from the Typing Result, including novel and failed alleles. This can be useful for downstream quality control of our sample. We will not go through the steps here, but a suggested analysis is to map the input sample reads to the alleles and then investigate the mapping of failed alleles, to evaluate whether the failed call is a true negative, or whether tool settings should be changed or additional references added to the MLST scheme, etc.

### Add sequence types to the ad hoc MLST scheme

After typing a sample, it is possible to add the result to an MLST scheme, which will, e.g., enable visual inspection of the relationship between samples.

We will now add the typing results to our ad hoc MLST scheme:

1. Launch **Add Typing Results to MLST Scheme** () using **Quick Launch** ()
2. In the first wizard step, "Select MLST Typing Results", select the seven successful typing results as input to the tool (leave out SRR16508428) (figure 17).

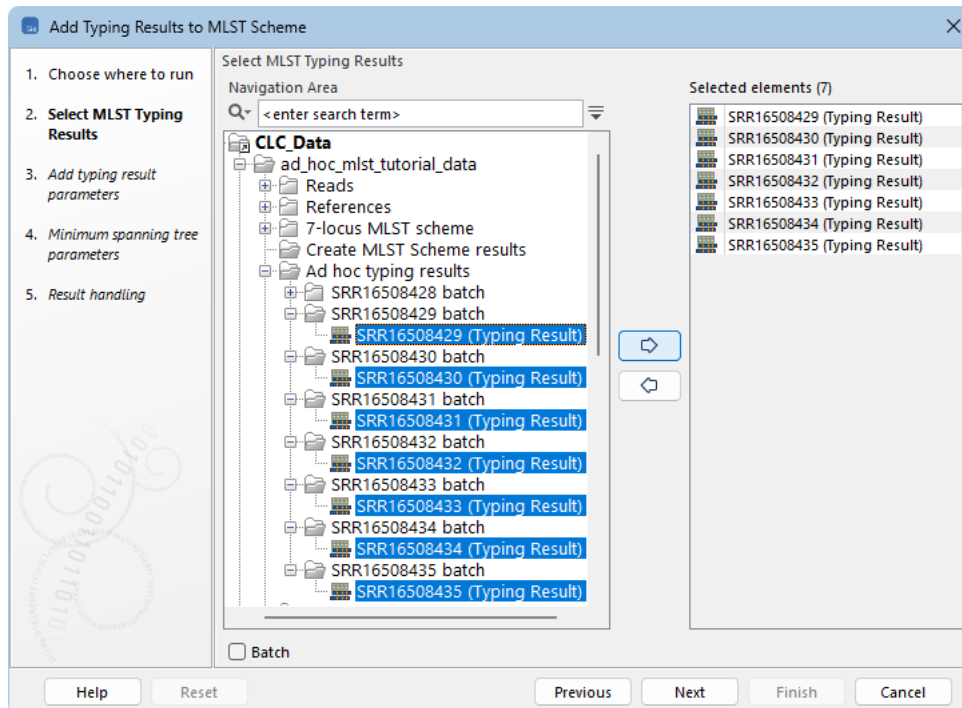
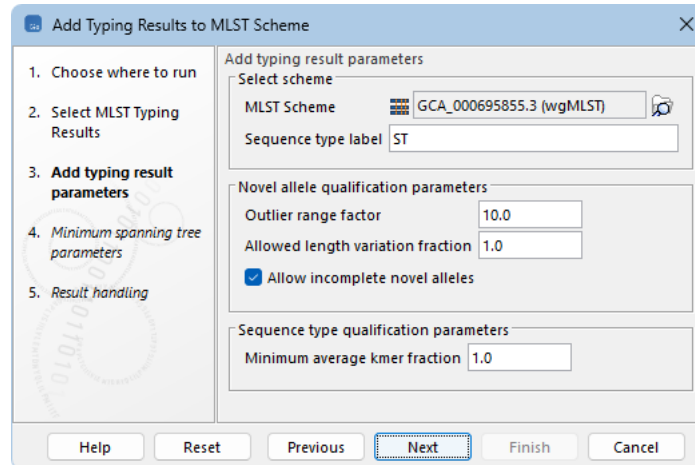


Figure 17: Select the seven successful typing results as input.

- In the next step, "Add typing result parameters", select the ad hoc MLST scheme we created earlier from the "Create MLST Scheme results" subfolder (figure 18).

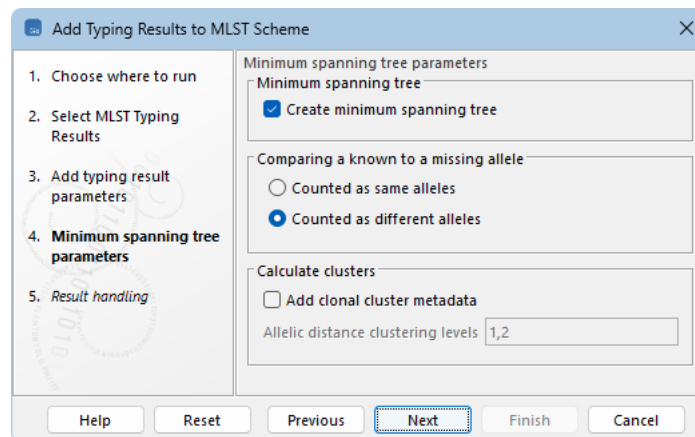
Set **Outlier range factor** = 10 and **Allowed length variation fraction** = 1 to allow all novel alleles to be included.



The screenshot shows the 'Add Typing Results to MLST Scheme' dialog box. The left sidebar lists five steps: 1. Choose where to run, 2. Select MLST Typing Results, 3. Add typing result parameters (highlighted), 4. Minimum spanning tree parameters, and 5. Result handling. The main panel is titled 'Add typing result parameters' and contains three sections: 'Select scheme' with 'MLST Scheme' set to 'GCA\_000695855.3 (wgMLST)' and 'Sequence type label' set to 'ST'; 'Novel allele qualification parameters' with 'Outlier range factor' set to 10.0, 'Allowed length variation fraction' set to 1.0, and 'Allow incomplete novel alleles' checked; and 'Sequence type qualification parameters' with 'Minimum average kmer fraction' set to 1.0. At the bottom are buttons for Help, Reset, Previous, Next (highlighted), Finish, and Cancel.

Figure 18: Updated "Add typing result parameters".

- In the next step, "Minimum spanning tree parameters", check **Create minimum spanning tree** (figure 19).



The screenshot shows the 'Add Typing Results to MLST Scheme' dialog box. The left sidebar lists five steps: 1. Choose where to run, 2. Select MLST Typing Results, 3. Add typing result parameters, 4. Minimum spanning tree parameters (highlighted), and 5. Result handling. The main panel is titled 'Minimum spanning tree parameters' and contains three sections: 'Minimum spanning tree' with 'Create minimum spanning tree' checked; 'Comparing a known to a missing allele' with 'Counted as different alleles' selected; and 'Calculate clusters' with 'Add clonal cluster metadata' unchecked and 'Allelic distance clustering levels' set to 1,2. At the bottom are buttons for Help, Reset, Previous, Next (highlighted), Finish, and Cancel.

Figure 19: "Create minimum spanning tree" must be checked.

- In the next step, "Result handling", select **Save**.
- In the last step, "Save location for new elements", click on the **New Folder** (📁) button to create a new subfolder in "ad\_hoc\_mlst\_tutorial\_data" and name it "Add Typing Results to MLST Scheme".

Choose to save the results in the new subfolder and click on **Finish**.

An updated MLST scheme and a report will be generated.

Note that if a publicly available wgMLST or cgMLST scheme already exists for the species, you can quickly assess the relationships between your samples using **Compare MLST Typing Results**. This tool can generate a minimum spanning tree and show allelic differences directly, without first needing to add typing results to an MLST scheme.

### Inspect the updated ad hoc MLST scheme

Open the updated **MLST scheme** (📄). The scheme contains three unique views, which can be opened by clicking on the icons in the bottom-left corner of the viewing area.

- The **Allele Table view** (📄) contains information about all loci and alleles in the scheme. Loci are listed in the upper table and all alleles for the selected locus or loci are listed in the lower table.
- The **Sequence Type Table view** (📄) provides an overview of the sequence types. There is a number of additional metadata columns, corresponding to the sample metadata available in the Metadata Table "A. baumannii isolates metadata", which was imported together with the samples. These columns are deselected by default, but can be selected from the **Side Panel** to the right of the view.
- The **Minimum Spanning Tree view** (🌳) is useful for visualizing relationships between strains or isolates. You can use the "Metadata" palette in the **Side Panel** to color the nodes by any of the metadata categories.

Open the Minimum Spanning Tree view (🌳). Then press **Ctrl** (⌘ on Mac) while clicking on the Sequence Type Table view icon (📄) to open both in split view (figure 20).

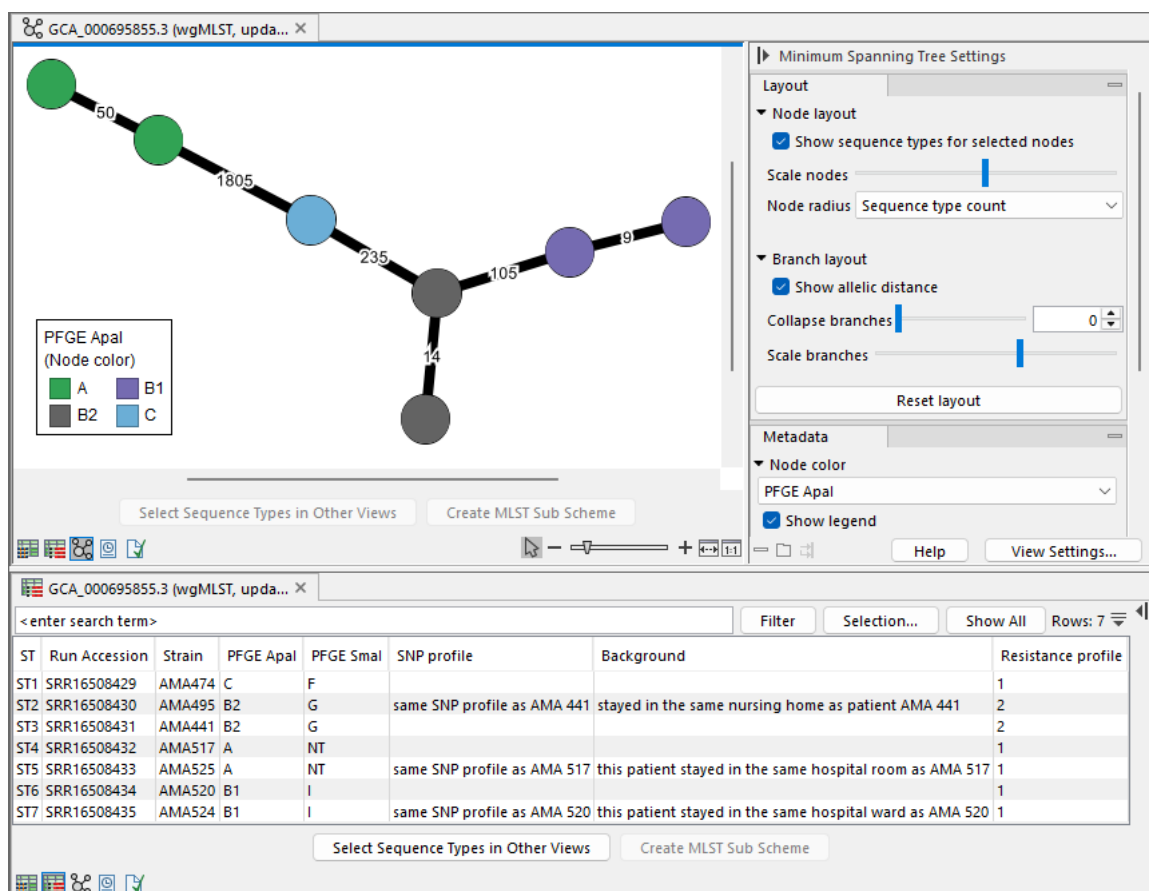



Figure 20: The minimum spanning tree (top) and metadata table (bottom) opened in split view. The seven successfully typed samples cluster in line with the metadata from [Hammerum et al., 2015].

In the **Side Panel** of the Minimum Spanning Tree view (figure 20 top), check the **Show sequence types for selected nodes** and **Show allelic distance** options. You can click and drag the "Scale branches" slider to the right to emphasize the distance between nodes. Under "Node color", select **PFGE Apal** and check **Show legend**.





The minimum spanning tree shows the distances between the samples based on the ad hoc MLST scheme (figure 20 top). Here, we can see that the shortest distances are between samples of patients that stayed in close proximity to each other according to the metadata. The clustering also correlates with the PFGE types and the SNP profiles (figure 20 bottom).

Using this updated MLST scheme for typing of new samples will make it easy to detect whether new samples are the same sequence type as existing samples. However, remember that this scheme is an ad hoc scheme, which is meant for analysis of samples that are thought to be related to the same outbreak.

If you are instead interested in creating a database of multiple outbreaks across larger areas (like multiple hospitals) or over longer periods of time, it is recommended to construct a more complete MLST scheme based on many references or use one of the well-established MLST schemes publicly available, e.g., by downloading them using [Download MLST Scheme](#) .

### Create an MLST scheme using a workflow

In the above sections, we stepped through each of the tools needed to create an ad hoc MLST scheme and add sequence types to it. This set of steps can be executed more efficiently by using the template workflow [Create MLST Scheme with Sequence Types](#), which can be found in the Toolbox:

**Workflows** | **Template Workflows**  | **Microbial Workflows**  | **Typing and Epidemiology**  | **Create MLST Scheme with Sequence Types** 

The template workflow is intended primarily for constructing a more complete MLST scheme. However, it can also support the creation of an ad hoc MLST scheme if the best reference(s) are already known.

If you want to inspect or edit the contents of the template workflow, you can [open a copy of the workflow](#) from the Toolbox.

## Bibliography

[Hammerum et al., 2015] Hammerum, A. M., Hansen, F., Skov, M. N., Stegger, M., Andersen, P. S., Holm, A., Jakobsen, L., and Justesen, U. S. (2015). Investigation of a possible outbreak of carbapenem-resistant *Acinetobacter baumannii* in odense, denmark using pfge, mlst and whole-genome-based snps. *Journal of Antimicrobial Chemotherapy*, 70(7):1965–1968.