



Tutorial

De Novo Assembly and BLAST

March 2, 2017

— Sample to Insight —

De Novo Assembly and BLAST

This tutorial takes you through some of the tools for a typical *de novo* sequencing workflow with a data set from a high-throughput sequencing machine. Here, we *de novo* assemble some reads, and then search a database at the NCBI with some of the contigs produced.

For this tutorial, we use an *E. coli* data set consisting just over 400,000 reads from a 454 sequencer.

Importing the data If you don't already have this sample data set:

1. Download the data set from our web site: http://resources.qiagenbioinformatics.com/testdata/raw_data/454.zip.
2. Unzip the file somewhere on your computer (on the Desktop for example).
3. Start the *CLC Genomics Workbench* and import the data:

File | Import (📁) | Roche 454

This will bring up the dialog shown in figure 1

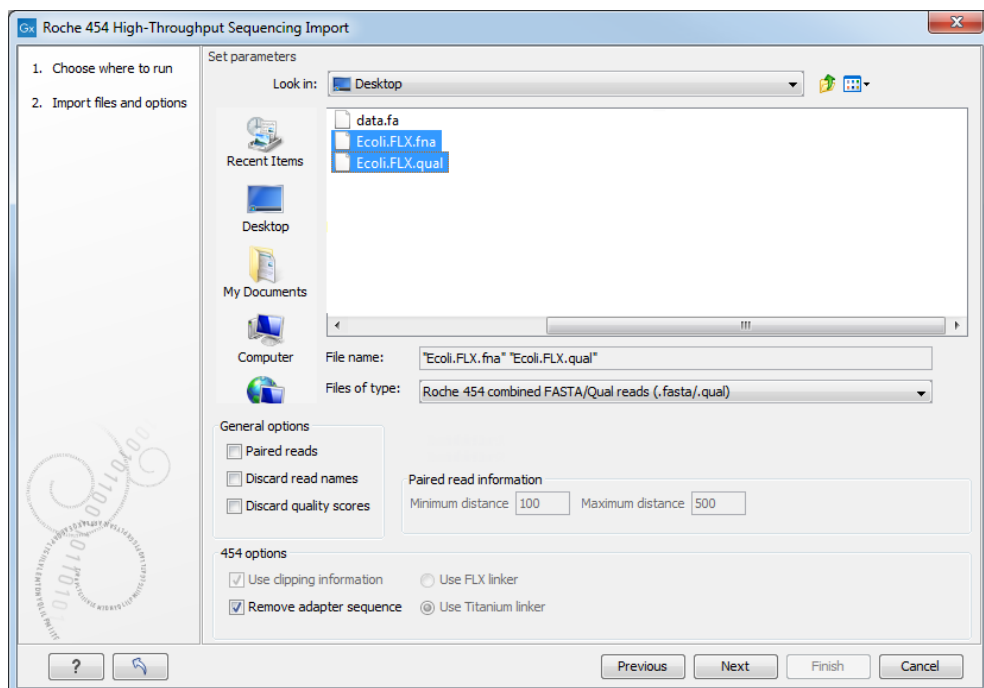


Figure 1: Choosing the file you wish to import.

4. Select the `Ecoli.FLX.fna` and `Ecoli.FLX.qual` files that come from the downloaded zip file. Make sure the **Remove adapter sequence** checkbox is checked and that the **Paired reads** checkbox is NOT checked. The option to discard read names is not significant in this context because of the relatively small amount of reads. So you can leave this checked, or unchecked, as you like. Click on the button labeled **Next**
5. Choose to **Save** the file, select the folder you wish to save to (you can create a new folder for this tutorial for example) and click on the button labeled **Finish**.

After a short while, the reads will be imported.

Assembly

The reads we are using in this tutorial are on average around 235 bp long. By doing a *de novo* assembly of the reads, we are trying to create longer, contiguous sequences from these relatively short sequences.

We recommend that you refer to the CLC Genomics Workbench manual for information on how the *de novo* assembly tool works, and the meaning of the parameters you have control over.

There are two possible types of output you can generate from the *de novo* assembly tool in the *CLC Genomics Workbench*: simple contigs, which is a sequence list of the contigs generated, and mapping objects, where a read mapping is carried out after the initial assembly, to map your reads back to the contigs created.

If you plan to do variant detection later, then generating mapping output can make sense. Also, there is an option when setting up your assembly, allowing you to update your contigs based on the mapping. If you enable this option, then areas of your assembled contigs that no reads map back to will be cut out of the final contigs. The idea there is that if no reads map back to a region of a contig, then there is no evidence in the data that that region exists. If this is what you want to do, then choosing to generate mapping output makes sense. You can, of course, always retrieve the contig sequences from mapping objects afterwards.

Sequence lists, like a list of contigs, can be used for different types of tasks, such as searching for patterns, or BLAST searching. As this is what we will proceed to do in this tutorial, we will choose to generate a sequence list of contig sequences in this tutorial. There is an optional section at the end of the tutorial about generating mapping output instead, and then extracting contig sequences from that for downstream use.

1. Start up a *de novo* assembly analysis by going to:

Toolbox | De Novo Sequencing (📁) | De Novo Assembly (🔧)

This shows the dialog in figure 2.

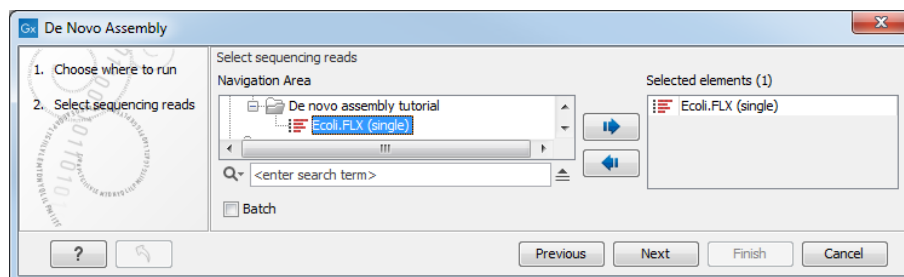


Figure 2: Select sequence list containing the reads.

2. Select the `Ecoli.FLX` (📁) sequence list and add it to the panel to the right. Click on the button labeled **Next**.
3. You now have the opportunity to set some parameters for your assembly. Here, we will accept the defaults, as shown in figure 3. Click on the button labeled **Next**.

You can directly access the manual to find out more about these parameters by clicking on the **Help** (?) button in the far left hand, bottom edge of the wizard window.

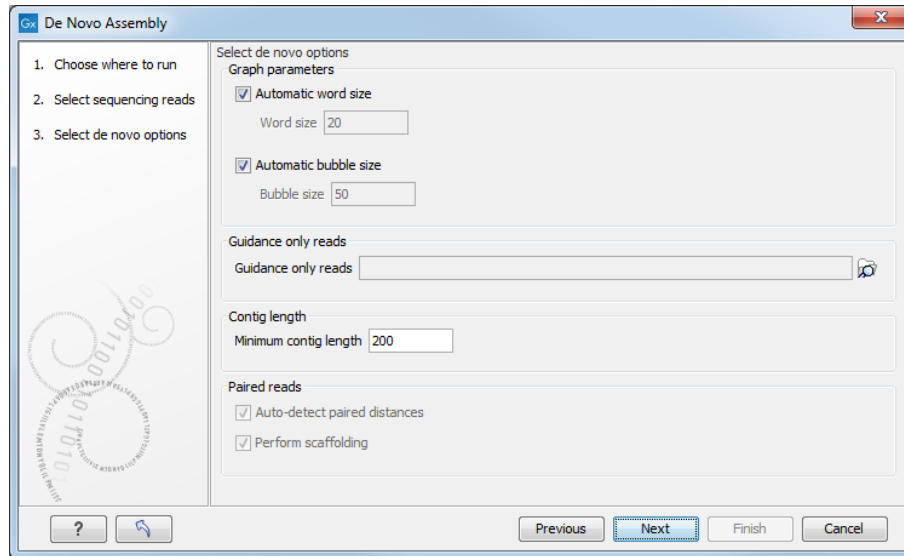


Figure 3: De novo assembly parameter options.

- On this next page of the wizard, you get to choose the type of output you will generate. As mentioned earlier, we will choose to create a **simple contig output**, i.e., a sequence list of the contig sequences assembled during this job. Check the option **Create simple contig sequences (fast)**, as shown in figure 4. Note that this causes all other options on this wizard page to be grayed out as they pertain to mapping reads back to the contigs. Click on the button labeled **Next**.

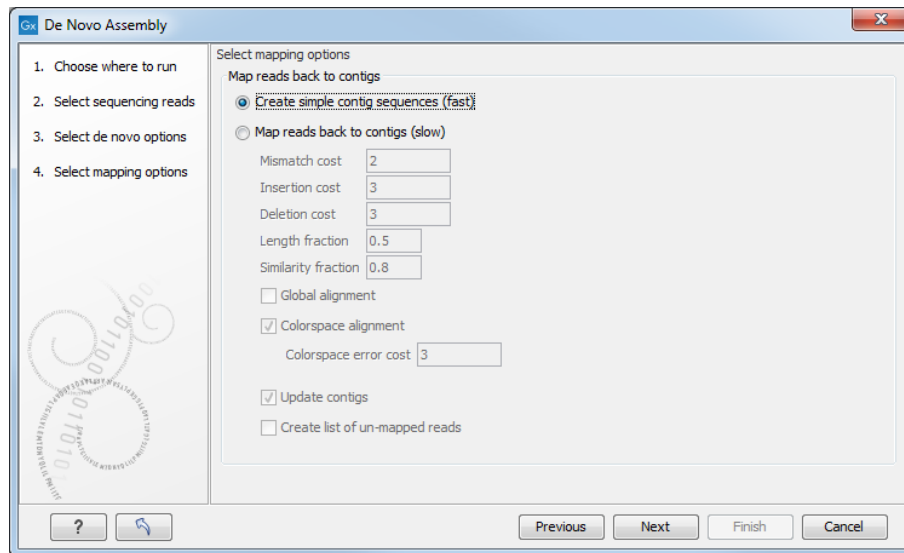


Figure 4: De novo assembly output types. Here we choose to create simple contigs. Note that if you choose to map reads back to the contigs, a read mapping job will be launched directly after the assembly task itself is done, and the outputs generated would be mapping objects.

- Choose to create a report and to Save the results. Select the location where you want to save the assembly output (your tutorial folder for example) and click on **Finish**.

The assembly should be quite quick, but this depends on how fast your computer is. You can follow the progress in the Processes panel of the workbench. Note that the progress bar does

not progress smoothly. It gives an indicator of the stage of the assembly more than an idea of the time to finish. For large assemblies, you may notice that the progress bar may spend a lot of time at a certain percentage. This generally means it is working on a certain stage of the assembly. You can see a brief text description of the stage it is working on written in brackets above the progress bar (see figure 5).

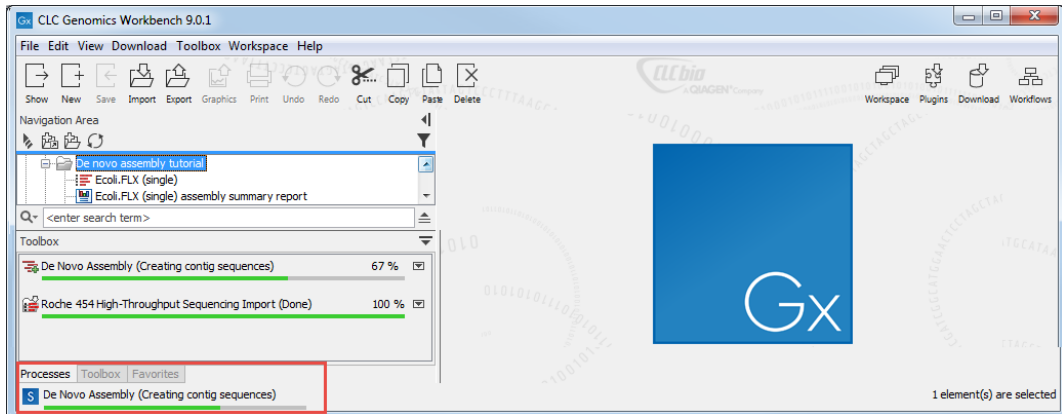


Figure 5: You can monitor the progress of an assembly using the Processes tab of the Toolbox in the workbench.

Investigate the results

The outputs of the assembly are a sequence list and a report. If you have not changed the name of the sample data, the contig list will be called **Ecoli.FLX (single) contig list**. The report will be called Summary de novo report.

1. Double click on the report object to open it. You will find general summary information in it such as the nucleotide distribution, information on contig lengths, as well as the N25, N50 and N75 values.
2. Double click on the sequence list you have created. This opens it in the View Area.
3. Click on the **table view icon** (📊) at the bottom of the workspace. You will now see a table view. You can see how many contigs were generated by how many rows are reported in the table.
4. Click on the size column heading to sort the table according to the length of the contigs. Highlight the top five rows - here, the longest contigs - as shown in figure 6.
5. Click on the button labeled **Create a New Sequence List**. This new sequence list is in a different tab in the View Area of the workbench.
6. Save this new list. You can do this in a number of ways. For example, you could:
 - Right click on the tab at the top of the new view, and choose **Save As...** from the menu that appears, or
 - Move the cursor over the tab at the top of the new view, press the left mouse button. Keeping the mouse button down, drag the tab over into a folder in the **Navigation Area** of the Workbench. The data will then be saved into that folder.

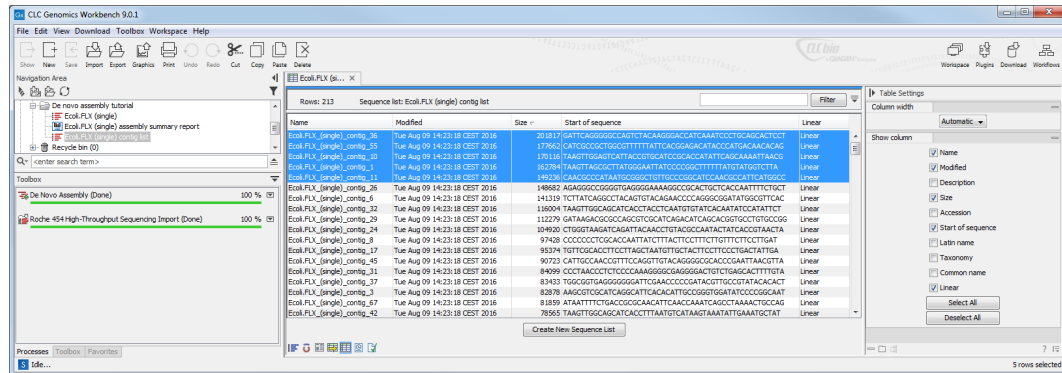


Figure 6: The table view of the sequence list, with the top five (longest) contigs highlighted.

BLAST some contigs against nr

One way to investigate the contigs would be to see if there are any similar sequences in public databases. This can be done using the BLAST program suite. In this tutorial we will run a BLAST search of the five contigs we saved into a new list earlier to search for similar sequences in a database at the NCBI. Note that you need to be connected to the internet to do this part of the tutorial.

1. Start up the BLAST tool in the Workbench by going to:
 - Toolbox | BLAST** (📁) | **BLAST at NCBI** (🌐)
2. Select the sequence list subset with the 5 contig sequences that you just saved. Click on the button labeled **Next**.
3. Choose **blastn** as the program type to run, and the **Nucleotide collection (nr)** as the database you wish to search. Click on the button labeled **Next**.
4. On this next page of the wizard, you can select the parameters you wish the blast search to be run with. There is some explanation of what these mean in the manual. You can access this by clicking on the **Help (?)** button in the far left hand, bottom edge of the wizard window.
5. Change the settings of the parameters so they match those in figure 7. Here, we set a longer word size, and a much smaller e-value than the default. We have also limited the maximum number of matches returned to us to 10. The sequences we have sent to the NCBI are quite long, so the search will take a little while.
6. Choose to **Save** your results, click on **Next**, choose a place to save your results to and then click on the button **Finish**.

Viewing the BLAST results

1. When the Blast search is finished, open the resulting data object. It should be called **Multi BLAST (5 sequences)**.
 You searched here with five query sequences. The table you just opened contains a brief overview of the results of your search for each of these query sequences.

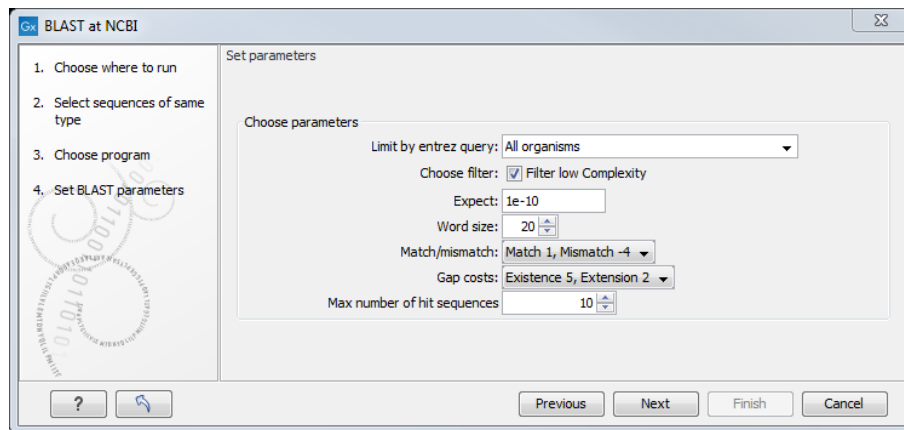


Figure 7: Setting parameters for a BLAST search.

- Look to the right of the table, in the View settings panel. You have many choices as to which columns you wish to see in this table. In figure 8, the Description column for the hit with the lowest e-value has been selected for display, in addition to the options checked by default. Note that you can choose to display your top hit according to e-value and/or according to identity.

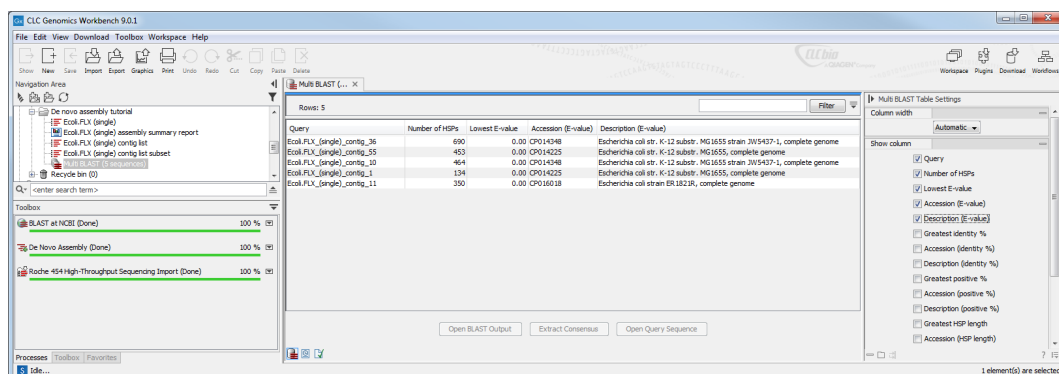


Figure 8: The overview BLAST result table.

For each contig, you can see the description of the top hit which is *E. coli* in each case, although of different strains.

You can, of course, export this summary table if you wish. Popular export formats for this sort of data would be excel, or comma separated values (csv).

- Double click on one of the rows of the summary blast table. This opens an individual BLAST result so you can investigate it in more detail.

By default, a visual depiction of the blast results will be shown. You can also click on the **table view icon** (📄) at the bottom of the workspace to see a table of results. With the table, one can filter and sort on different values. Like any table, you can export the data, with popular formats being excel, or comma separated values (csv).

If you are used to standard, text blast format, you may wish to click on the **Text Contents** (📄) icon at the bottom of the workspace.

Choosing mapping as the output type from an assembly

This is an optional part of the tutorial. Here, we re-run the de novo assembly, but this time choosing to map the reads back to the contigs, and to update the reads based on the mappings.

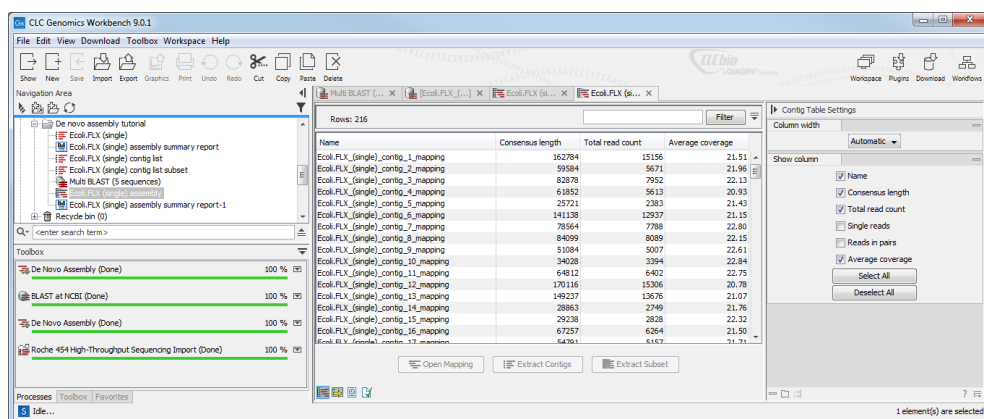
1. Start up a de novo assembly analysis by going to:

Toolbox | De Novo Sequencing (🧬) | De Novo Assembly (🧬)

2. Select the `Ecoli.FLX` (📄) sequence list as input by adding it to the panel to the right and click on **Next**.
3. You now have the opportunity to set some parameters for your assembly. Here, we will accept the defaults, as shown in figure 3. Click on **Next**.
4. On this next page of the wizard, you get to choose the type of output you will generate. Here, choose **Map reads back to contigs (slow)**. Leave all the mapping options as the defaults. Please refer to the manual, for example, by clicking on the **Help (?)** button to find out more about these options. Also add a check in the box beside **Update contigs** before clicking on **Next**.
5. Choose to **Create a report** and to **Save** the results in your tutorial folder.

The outputs of the assembly are a summary table of the mappings and a report. If you have not changed the name of the sample data, the mapping output will be called **Ecoli.FLX (single) assembly**. The report will be called **Ecoli.FLX (single) assembly summary report-1**. There are a couple of extra graphs in this report compared to the one you generated earlier, where you requested only the contigs themselves as output.

1. Open the mapping output by double clicking on it in the **Navigation Area**. You should see a view like that shown in figure 9.



Name	Consensus length	Total read count	Average coverage
Ecoli.FLX (single)_contig_1_mapping	162784	15156	21.51
Ecoli.FLX (single)_contig_2_mapping	99584	5671	21.96
Ecoli.FLX (single)_contig_3_mapping	62870	7952	22.13
Ecoli.FLX (single)_contig_4_mapping	61852	5613	20.93
Ecoli.FLX (single)_contig_5_mapping	25721	2383	21.43
Ecoli.FLX (single)_contig_6_mapping	141138	12937	21.15
Ecoli.FLX (single)_contig_7_mapping	78564	2768	22.80
Ecoli.FLX (single)_contig_8_mapping	84099	8089	22.15
Ecoli.FLX (single)_contig_9_mapping	51084	5007	22.61
Ecoli.FLX (single)_contig_10_mapping	34020	3394	22.84
Ecoli.FLX (single)_contig_11_mapping	64812	6402	22.75
Ecoli.FLX (single)_contig_12_mapping	170116	15306	20.78
Ecoli.FLX (single)_contig_13_mapping	149237	13676	21.07
Ecoli.FLX (single)_contig_14_mapping	28863	2749	21.76
Ecoli.FLX (single)_contig_15_mapping	29238	2828	22.32
Ecoli.FLX (single)_contig_16_mapping	67257	6264	21.50
Ecoli.FLX (single)_contig_17_mapping	64761	6182	21.71

Figure 9: A summary mapping table.

Each row in the table represents a contig. There should be around 200 contigs (the number can vary a little from version to version). Double-clicking a row will open the contig with the mapped reads.

2. Sort the table on Average coverage by clicking on the column heading. Sort them so that the contigs with the highest average coverage appear at the top of the table.

3. Double click on the top contig listed (i.e. the one with the highest coverage) to open it in the viewing area.
4. Maximize the view of this assembly by double clicking on the tab name in the Viewing area.
5. Feel free to investigate the effects of the different viewing options on how the mapping looks, and can be searched. Please also try zooming in at different levels.
6. When you have seen enough, restore the view again by double clicking on the tab name again.
7. Close the assembly tab - for example, by clicking on the small X at the right side of the tab.
8. Go back to viewing the mapping summary table.
9. Sort the mapping table on Consensus length.
10. Highlight the rows containing information about the five longest contigs.
11. Click on the button marked **Extract Contig** to get a sequence list with just these contigs. Choose to **Save** these.

Now you have a sequence list for your 5 longest contigs. From here, you could run a BLAST job, or any other task that requires a sequence list as input.