

BioSeqDB Help Document

April 27, 2021

Arnie Berg

Table of Contents

Background	2
File systems	3
Navigating file systems with the Explorer	4
About dialog	6
Login dialog	7
Functions	8
0. New database	8
1. Assemble	9
2. Insert	11
3. Extract	12
4. Remove	12
5. Delete database	13
6. Backup database	13
7. Restore database	14
Analysis Functions	15
0. BBMap	15
1. Build tree	16
2. Influenza A pipeline	19
3. Kraken2	21
4. Quast	22
5. Salmonella Serotyping	23
6. Search	24
7. VFabricate	25
Running in the background	26
Editing LIMS identifiers	28
For BioSeqDB administrators	29
FAQ	33
Appendix A: Deploying the BioSeqDB client	35

Background

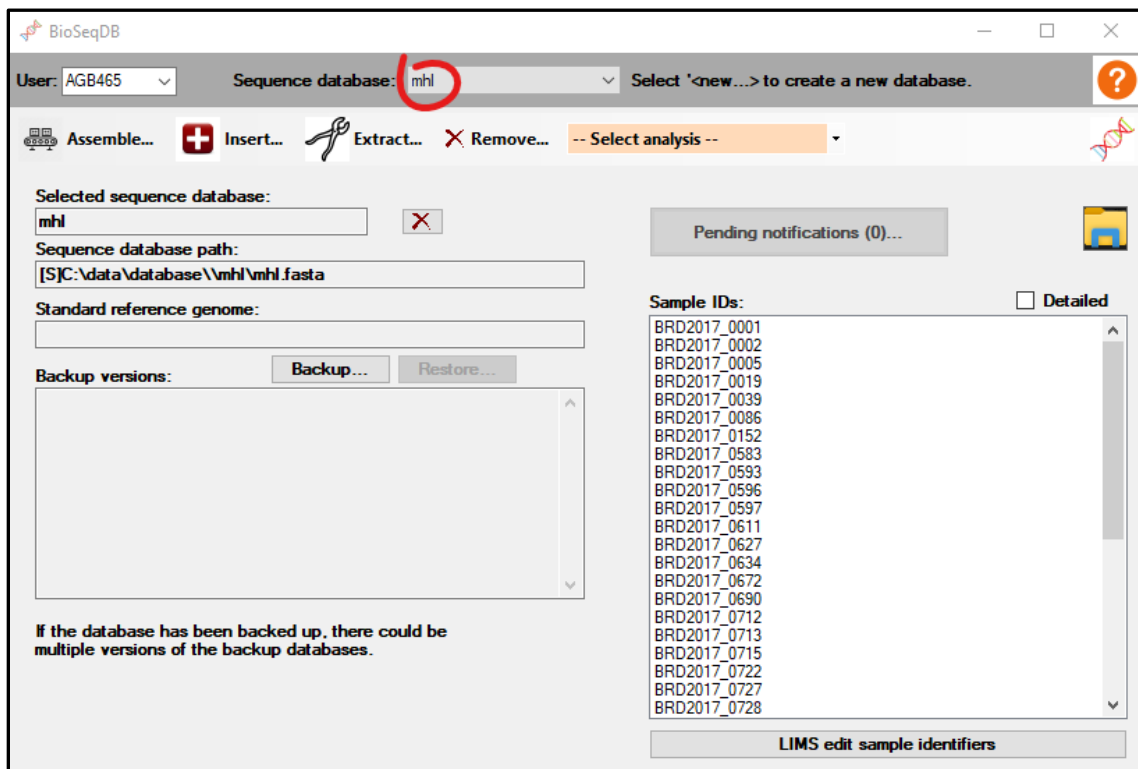
The purpose of the BioSeqDB program is to manage the process flow and information repository for sequence data of identified organisms and provide access to functional analysis applications that run in the Linux WSL environment.

BioSeqDB runs as a Windows program and appears on the task bar as follows:



The BioSeqDB program is basically a wrapper that invokes Linux commands that have been developed to run in the Windows Subsystem for Linux (WSL) environment. The advantage of running BioSeqDB instead of the individual Linux commands is that BioSeqDB prompts the user for all the relevant parameters and functions needed to accomplish the biologically important objectives and remembers previously selected options.

BioSeqDB maintains a separate database for each specified organism. If no existing database exists in which to insert a sequence for a different organism, a new database must be created. The BioSeqDB dialog always displays the currently selected sequence database, and optionally a standard reference genome associate with the database and a list of database backup versions. Databases may be selected from the 'Sequence database' dropdown. The main BioSeqDB dialog also displays the Sample IDs from the currently selected database and appears as follows:



BioSeqDB is designed as a multi-user system with clients accessing the BioSeqDB service (running on WIMMER) over the network. The databases reside on the server to be centrally available to all users.

Most of this document describes the functionality available through this dialog. A section near the end provides some details for administrators of BioSeqDB. This document is meant to assist users in running BioSeqDB. It does not describe how to interpret results of the functions invoked by BioSeqDB.

File systems

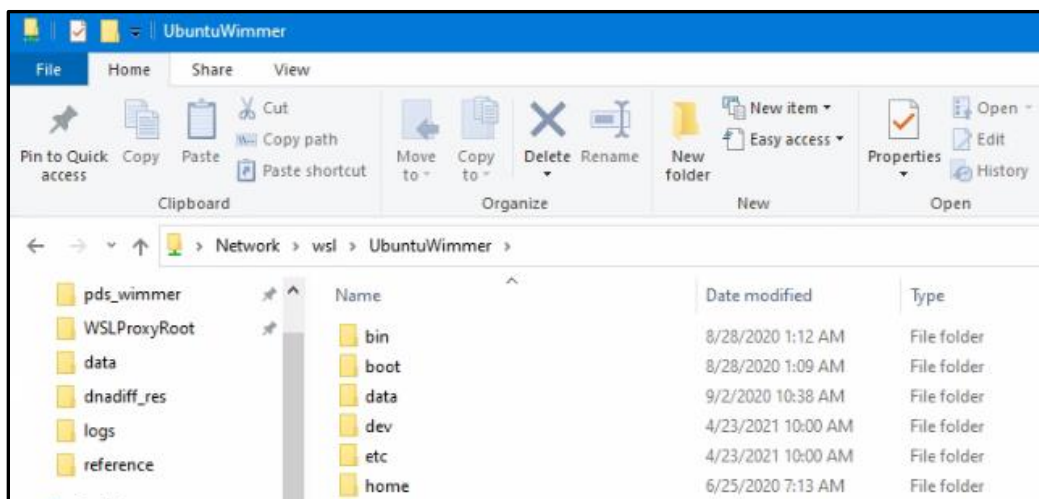
Through WSL, BioSeqDB provides complementary access to the Windows file system and the Linux file system. BioSeqDB also provides access to both the client computer file system and the server file system.

If you are familiar with Linux, by connecting to WIMMER via SplashTop and using Windows Terminal in Windows you can start up the UbuntuWimmer shell as Linux user pds_wimmer to run in the Linux environment. Those details are hidden by BioSeqDB, but it still is important to be aware of the presence of the Linux file system and how it relates to the Windows file system.

From Windows File Explorer on WIMMER, the Linux file system may be accessed by entering '\\wsl' in the address bar. This shows an instance of the Linux file system called UbuntuWimmer:



Double-click on UbuntuWimmer to open the top-level folder structure of the Linux file system:



If you are familiar with Windows, you should be aware of how Linux references folders in the Windows file system. This is important for BioSeqDB because although most functions run in Linux, most data are stored in the Windows file system. For example, all sequence data, both raw and assembled, are stored on the E: drive of WIMMER.


References to the C: and E: drive of WIMMER are represented as /mnt/c and /mnt/e respectively in Linux. For example, a file in the Temp folder on the C: drive might be C:\Temp\stats.txt in Windows, but in Linux would be referenced as /mnt/c/Temp/stats.txt. In the Linux reference, upper-case and lower-case is important. In Windows it is not. Also note the difference in the use of the backslash in the Windows file system and the forward-slash in the Linux file system.


The Explorer (described later) in BioSeqDB hides most of these details and translates folder references automatically in the background but understanding these differences between the file systems is important.

A key rule to keep in mind is that no path or file name may contain white space (like a space character).

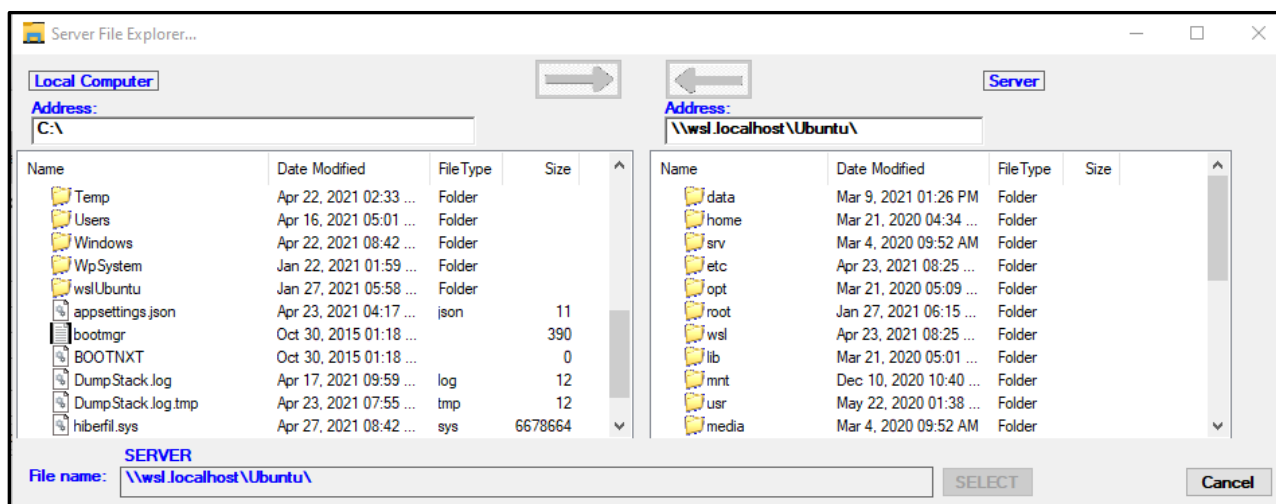
In BioSeqDB, there is also a distinction between files and folders stored on the server and those stored on the client. For example, the output results from running an analysis in BioSeqDB may be stored on either the server or the client computer. To distinguish which file system is intended, the path of the file or folder is prefixed in BioSeqDB with either a '[S]' if on the server, or a '[L]' if on the local computer.

Navigating file systems with the Explorer

The Explorer is frequently opened to locate data files and folders on the file system, whether in Windows, Linux, or on the local or server computer. Individual functions within BioSeqDB often provide ellipses to invoke the Explorer (). For general access, the Explorer can be opened by clicking on the

folder icon () on the right side of the main BioSeqDB dialog.

Since the user needs to have detailed access to the file system on the server, BioSeqDB has a unique file and folder Explorer to navigate both the local computer and the server file system. The local computer



is described on the left side of the dialog and the server is described on the right side. This layout should be familiar to users of Globus, FTP or various other remote connection tools.

Double-click on any folder name to drill down to the next level. At times, the Explorer dialog is opened from a function that is looking for a file and other times for a folder. The 'SELECT' command is disabled until a valid selection of a file or folder is made. When the 'SELECT' command is enabled and clicked, the dialog is closed with the currently selected path returned to the application.

Right-click on either the Local Computer or Server part of the Explorer window to display a context menu that allows you to change the view or delete files or folders.

Click on the 'Local Computer' label or 'Server' label at the top of the dialog to reset the explorer back to the contents of the root folder for the respective computer.

Note that this Explorer is capable of drilling deep into the WSL Linux file structure as easily as the Windows file structure. Also, if any mapped drives exist on either local computer or the server, that drive shows up as well.

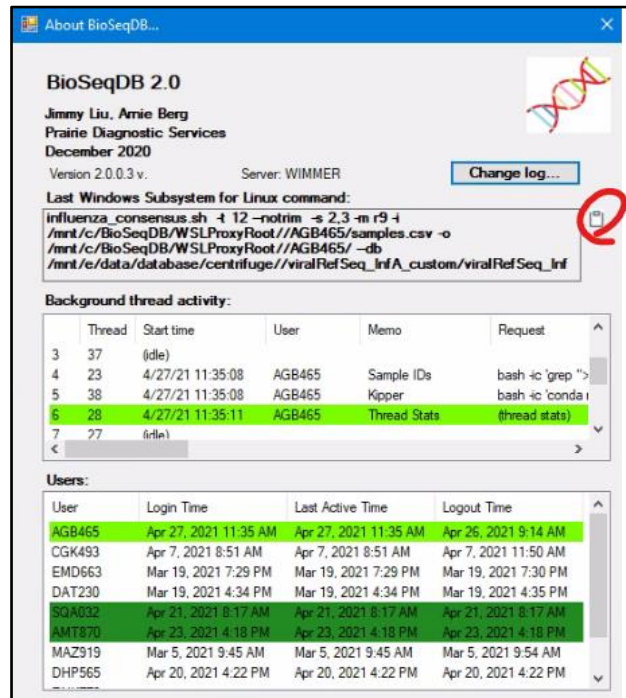
There are two arrows at the top of the dialog. They become enabled any time that a file transfer or folder transfer from the local computer to the server (or vice versa) is valid. This allows transfer of data between computers if needed. Multiple folders/files may be selected for transfer between computers and drag-and-drop with the mouse is also supported. Be aware though that if the analysis you are running gives you the option to specify a folder as output or a file as input *on your local computer*, BioSeqDB automatically and seamlessly looks after transferring those files between computers.

About dialog

The About dialog for BioSeqDB 2.0 can be opened by clicking on the DNA double-helix symbol on the main dialog. This dialog identifies the version of BioSeqDB along with development credits and the server to which the BioSeqDB is connected, but there are a couple of items of added value.

The 'Change log' command opens a list of changes made to BioSeqDB over time. This is helpful if the user wants to know what the most recent changes have been.

The 'Last Windows Subsystem for Linux command' lists the last command that was issued to WSL on behalf of the client. For testing and debugging purposes, this command may be copied to the clipboard by clicking on the highlighted clipboard icon.



BioSeqDB 2.0
Jimmy Liu, Amie Berg
Prairie Diagnostic Services
December 2020
Version 2.0.0.3 v. Server: WIMMER [Change log...](#)

Last Windows Subsystem for Linux command:
`influenza_consensus.sh -t 12 -notrim -s 2.3 -m r9 -i /mnt/c/BioSeqDB/WSLProxyRoot//AGB465/samples.csv -o /mnt/c/BioSeqDB/WSLProxyRoot//AGB465/-db /mnt/e/data/database/centrifuge//viralRefSeq_InfA_custom/viralRefSeq_Inf`

Background thread activity:

Thread	Start time	User	Memo	Request
3	37	(idle)		
4	23	4/27/21 11:35:08	AGB465	Sample IDs bash -c 'grep ">
5	38	4/27/21 11:35:08	AGB465	Kipper bash -c 'conda i
6	28	4/27/21 11:35:11	AGB465	Thread Stats (thread stats)
7	27	(idle)		

Users:

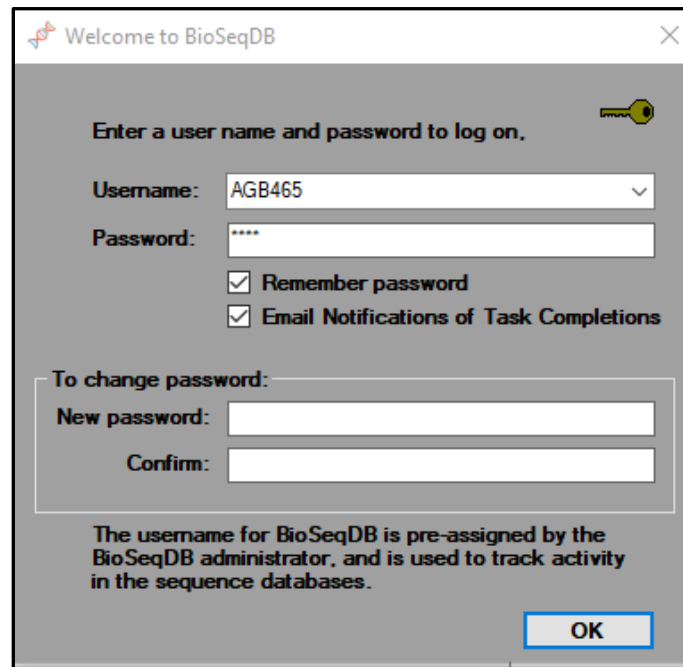
User	Login Time	Last Active Time	Logout Time
AGB465	Apr 27, 2021 11:35 AM	Apr 27, 2021 11:35 AM	Apr 26, 2021 9:14 AM
CGK493	Apr 7, 2021 8:51 AM	Apr 7, 2021 8:51 AM	Apr 7, 2021 11:50 AM
EMD663	Mar 19, 2021 7:29 PM	Mar 19, 2021 7:29 PM	Mar 19, 2021 7:30 PM
DAT230	Mar 19, 2021 4:34 PM	Mar 19, 2021 4:34 PM	Mar 19, 2021 4:35 PM
SQA032	Apr 21, 2021 8:17 AM	Apr 21, 2021 8:17 AM	Apr 21, 2021 8:17 AM
AMT870	Apr 23, 2021 4:18 PM	Apr 23, 2021 4:18 PM	Apr 23, 2021 4:18 PM
MAZ919	Mar 5, 2021 9:45 AM	Mar 5, 2021 9:45 AM	Mar 5, 2021 9:54 AM
DHP565	Apr 20, 2021 4:22 PM	Apr 20, 2021 4:22 PM	Apr 20, 2021 4:22 PM

Another bonus item is a snapshot of the background thread activity. All requests to the Linux system are handled by a component on the server called WSLProxy. At this time there can be up to eight simultaneous active threads. This does not limit the number of CPU cores that the analysis can use but refers only to the number of simultaneous activities that BioSeqDB through WSLProxy can manage. The snapshot provides a report of what request is running on each thread, along with the requesting user and the start time. The list can be refreshed by clicking on the DNA double-helix symbol in the top right-hand corner of the About dialog. The thread highlighted in green is the thread that is used to report on this list and is always present.

The Users list shows the status of current users on the system. Reported is login, active and logout times, with currently active users highlighted in green.

Login dialog

BioSeqDB authenticates each user of the application (for administrators, see the section below on 'User Management'). The username is selected from a dropdown list and the password is entered to be validated. If the password is valid, it can also be changed by supplying a new password in the 'New password' and 'Confirm' fields.



The image shows a 'Welcome to BioSeqDB' login dialog box. It has a title bar with a close button. The main area is grey and contains the following elements: a key icon, the instruction 'Enter a user name and password to log on.', a 'Username:' dropdown menu with 'AGB465' selected, a 'Password:' text field with four asterisks, two checked checkboxes labeled 'Remember password' and 'Email Notifications of Task Completions', a section titled 'To change password:' containing 'New password:' and 'Confirm:' text fields, and a paragraph of text at the bottom stating: 'The username for BioSeqDB is pre-assigned by the BioSeqDB administrator, and is used to track activity in the sequence databases.' An 'OK' button is located at the bottom right.

Welcome to BioSeqDB

Enter a user name and password to log on.

Username: AGB465

Password: ****

☒ Remember password

☒ Email Notifications of Task Completions

To change password:

New password:

Confirm:

The username for BioSeqDB is pre-assigned by the BioSeqDB administrator, and is used to track activity in the sequence databases.

OK

In addition, there is a checkbox to 'Remember password'. This is useful if BioSeqDB is being used on your personal computer or laptop. The username is automatically remembered. The other checkbox indicates whether to send out 'Email Notifications of Task Completions'. Some operations in BioSeqDB may be of long enough duration that you may leave it running while attending to other matters. In that case, if 'Email Notifications' is checked, you are notified when the long running task is complete.

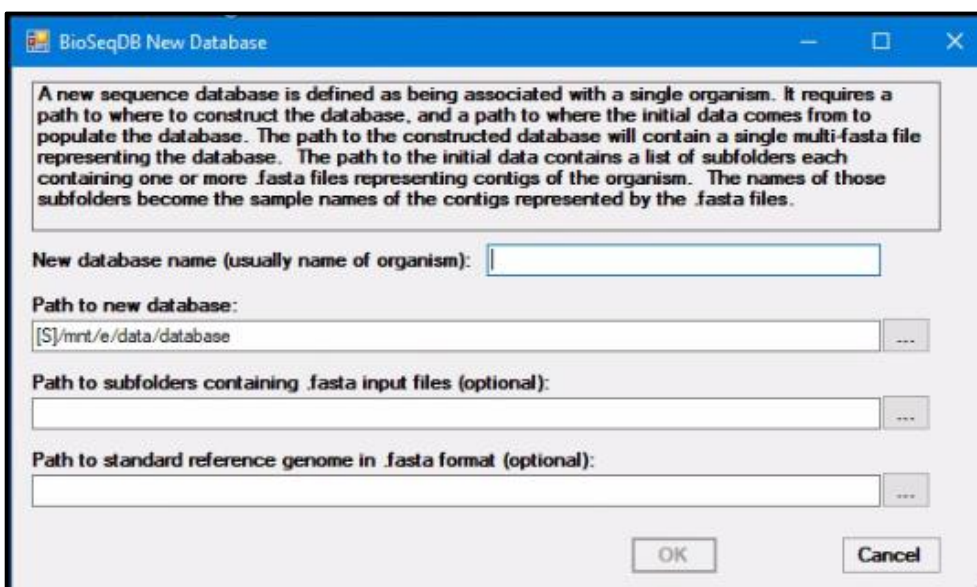
Functions

The menu strip at the top of the BioSeqDB dialog has the list of functions available, including Assemble, Insert, Extract, and Remove. From the main dialog, command buttons are available to Delete, Backup and Restore the selected database. From the '—Select analysis —' dropdown, analysis functions are selected to perform specialized analyses, such as BBMap, Build tree, Influenza A pipeline, Kraken2, Quast, Salmonella, Search and VFabricate. These analyses are described separately below.

A new database may be created by selecting '<new...>' from the 'Sequence database:' dropdown and supplying the sequence database setup information.

0. New database

A new database requires a database name, which is usually the name of the organism or its abbreviation. The second parameter is the path to the new database. By convention, sequence databases are stored in the E:\data\database folder on WIMMER. The third parameter is the path to the input sequence(s) representing the initial content of the database. This data may be stored anywhere, but it is important to understand the structure of this data.



The screenshot shows the 'BioSeqDB New Database' dialog box. It contains a text area with instructions: 'A new sequence database is defined as being associated with a single organism. It requires a path to where to construct the database, and a path to where the initial data comes from to populate the database. The path to the constructed database will contain a single multi-fasta file representing the database. The path to the initial data contains a list of subfolders each containing one or more .fasta files representing contigs of the organism. The names of those subfolders become the sample names of the contigs represented by the .fasta files.' Below this are four input fields: 'New database name (usually name of organism):', 'Path to new database:' (with the text '[S]/mnt/e/data/database'), 'Path to subfolders containing .fasta input files (optional):', and 'Path to standard reference genome in .fasta format (optional):'. Each field has a browse button (three dots). At the bottom are 'OK' and 'Cancel' buttons.

Specify the path containing one or more subfolders, where each subfolder contains one or more .fasta files representing contigs or consensus sequences of the organism. The names of the subfolders are the sample IDs that are used to name the sample in the new database. Unless this structure is set up ahead of creating the new database, the setup of the new database will not succeed. To create an empty sequence database, leave this path empty.

Optionally a standard reference sequence genome may be specified to be associated with the sequence database. This reference genome is useful when running the Build tree function for this database.

1. Assemble

To create new data to add to a sequence database, raw sequence data must be assembled to create contig data in .fasta files. The Assemble function performs this step.

The Assemble function can assemble multiple sets of sequence data at one time. In fact, it is desirable to assemble multiple sequence data at once to take advantage of the multi-tasking capability of WIMMER. Use the Sample Picker to select the folder containing the .fastq files to include in the Assemble step. Only samples that are checked in the Samples list are included in the Assemble function. Any unchecked samples can be deleted with the 'Delete unchecked' command button. The sample name is derived from the immediate parent folder of the .fastq data. For example, in the selected sample below, the sample name of the selected sample is 'subset1_2040649_1'.

The screenshot shows the 'BioSeqDB Assemble Samples' dialog box. At the top, a text box explains: 'Assemble samples using NextFlow. Select folders containing .fastq files to assemble with the Sample Picker. Final assemblies with analyses are stored under the staging folder using the sample name as the subfolder name. Optionally specify the maximum number of fastq files to pick from each sample.'

Samples to assemble: A list of sample folders is shown. The first item, 'E:\data\Anatoly\PDS2040649\subset1_2040649_1\' is selected with a checkmark. To the right of this list is a 'Max # fastq files:' input field. To the right of the list are two buttons: 'Sample picker...' and 'Delete unchecked'.

Select assembler: Two radio buttons are present: 'Rapid Assembler' and 'Flye'. Below them, under 'Analyses to perform:', are four checkboxes: 'Kraken2' (checked), 'BBmap' (checked), 'Quast' (checked), and 'VFabricate' (unchecked). Below this is a text field for 'Gene cross-reference configuration table (for VFabricate):' with the value '[S]E:\data\reference\VFgeneXRef.csv'. A checkbox 'Use fast polish (skip Medaka step)' is also present.

Trinity (viral): This section is selected with a radio button. It contains a 'Reference genomes:' section with a 'Virus reference:' text field containing '[S]\wsl\$\UbuntuWimmer\home\pds_wimmer\JimmyFiles\Influenza_consensu' and a 'Host reference (optional):' text field.

At the bottom, there is a 'Memo (to help identify task in notifications):' text area, and 'OK' and 'Cancel' buttons.

Two bacterial assemblers are available, either Flye or Rapid Assembler. Although they produce similar results, sometimes one will fail to successfully complete the assembly and the other will succeed. There is also an optional Medaka step which further polishes the data but takes an increased amount of time. Any of four analyses to perform can be selected, Kraken2, BBmap

Quast and VFabricate. If VFabricate is selected, the gene cross-reference configuration table can also be defined and accessed for editing.

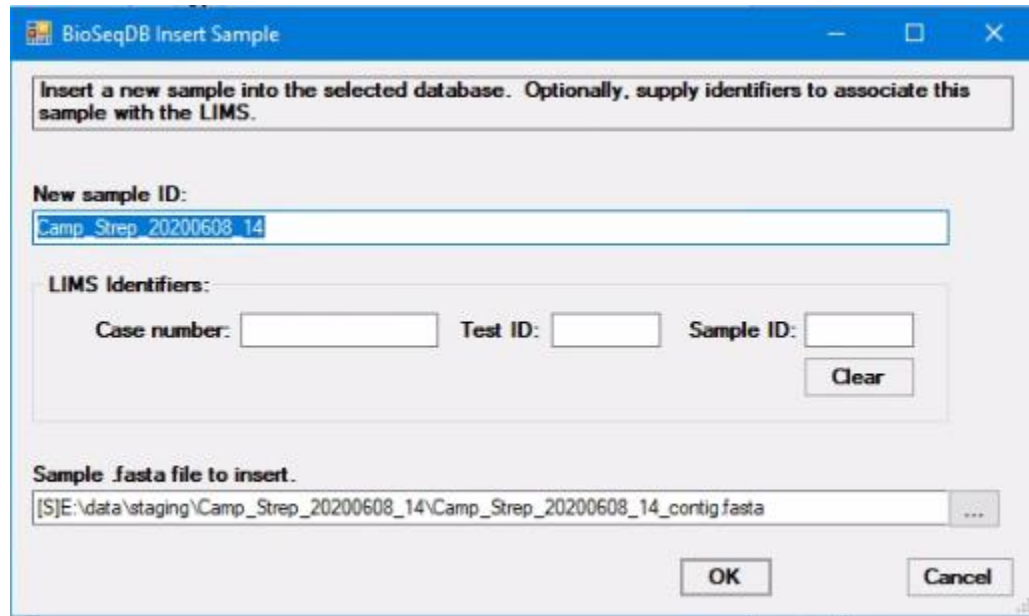
A viral assembler based on Trinity is also available. The assembly is based on the reference genome for the virus. If the sequence data contain host sequences, a host reference genome may optionally be specified to remove the host DNA data. Because this assembly uses porechop to demultiplex and polish the data, this assembly can run for hours, depending on the amount of sequence data. This assembly workflow requires that the viral data be barcoded.

The results of the Assemble function, if successful, are stored in the E:\data\staging folder, with a subfolder created named after the sample name of the data assembled. A fasta contig file is created that can then be inserted into a sequence database for that organism.

The Assemble function is one of the functions that can be quite time-consuming. As such, it is scheduled to run in the background and its status is referred to as 'Pending' until it completes and becomes 'Ready'. The 'Pending Notifications' dialog provides the ability to manage this process. The details about scheduled functions are provided below in 'Running in the background'.

The optional 'Memo' field at the bottom of the Assemble dialog can be used to associate details relevant to the assembly to better track the task in the background. There is also a 'Max # fastq files' field where the number of fastq files to be read can be limited. This allows you to process a subset of the full set of fastq files without having to create a set of subset files.

2. Insert



The image shows a Windows-style dialog box titled "BioSeqDB Insert Sample". It contains the following elements:

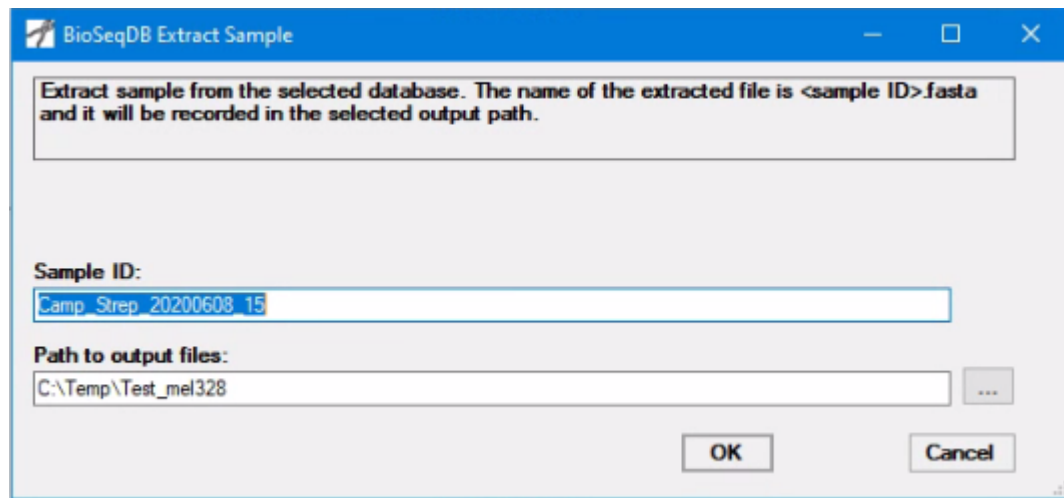
- A message box at the top: "Insert a new sample into the selected database. Optionally, supply identifiers to associate this sample with the LIMS."
- A section labeled "New sample ID:" with a text input field containing "Camp_Strep_20200608_14".
- A section labeled "LIMS Identifiers:" containing three input fields: "Case number:", "Test ID:", and "Sample ID:". There is a "Clear" button to the right of these fields.
- A section labeled "Sample .fasta file to insert." with a text input field containing the file path "[S]E:\data\staging\Camp_Strep_20200608_14\Camp_Strep_20200608_14_contig.fasta" and a browse button (three dots) to its right.
- At the bottom right, there are "OK" and "Cancel" buttons.

The Insert function inserts the contents of a .fasta consensus sequence or contig file into the current database. You are prompted to create a new sample ID and select the .fasta file containing the contig(s).

If the sample ID you enter already exists in the database, you are prompted as to whether you want to replace the existing data associated with the sample ID in the database.

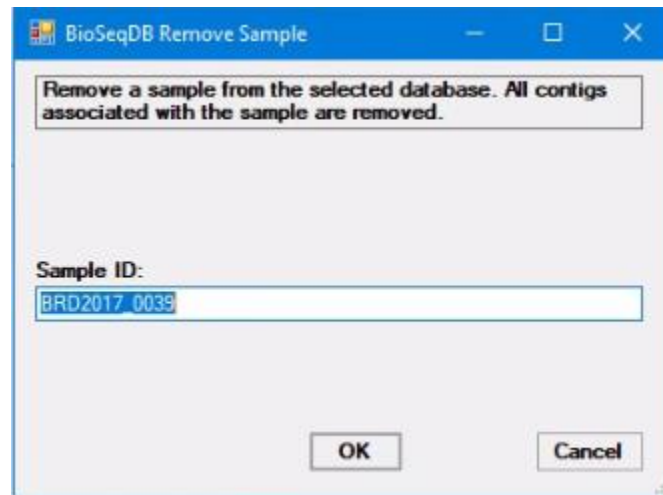
Optionally, to connect the sequence sample with the sample in the Casebook LIMS, a case number, LIMS test ID and LIMS sample ID may be specified. These LIMS identifiers may also be edited by clicking on the 'LIMS edit sample identifiers' command button on the main dialog to open a dialog to edit values for any sequence sample in the currently selected sequence database. For more details, see the section below entitled 'Editing LIMS identifiers'.

3. Extract



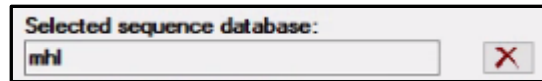
The Extract function is used to extract a sequence from the currently selected database. You are prompted for a Sample ID and a path where the extracted sequence can be recorded. Since the sample is selected from the currently selected database, it is enough just to select the sample ID from the sample ID list on the main dialog.

4. Remove



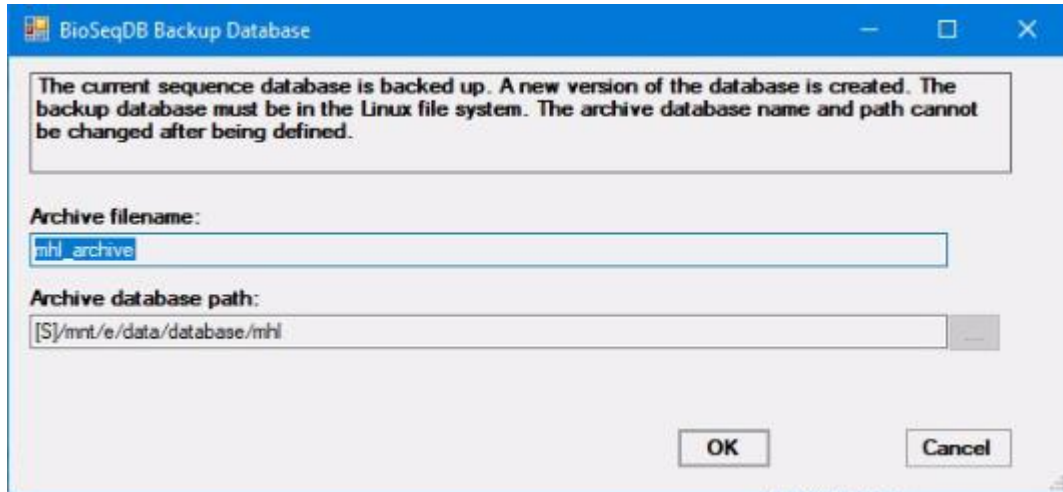
As with the Extract function, the Remove function uses the specified Sample ID, either entered from the keyboard or selected from the sample ID list. The specified Sample ID is removed from the sequence database if it exists.

5. Delete database



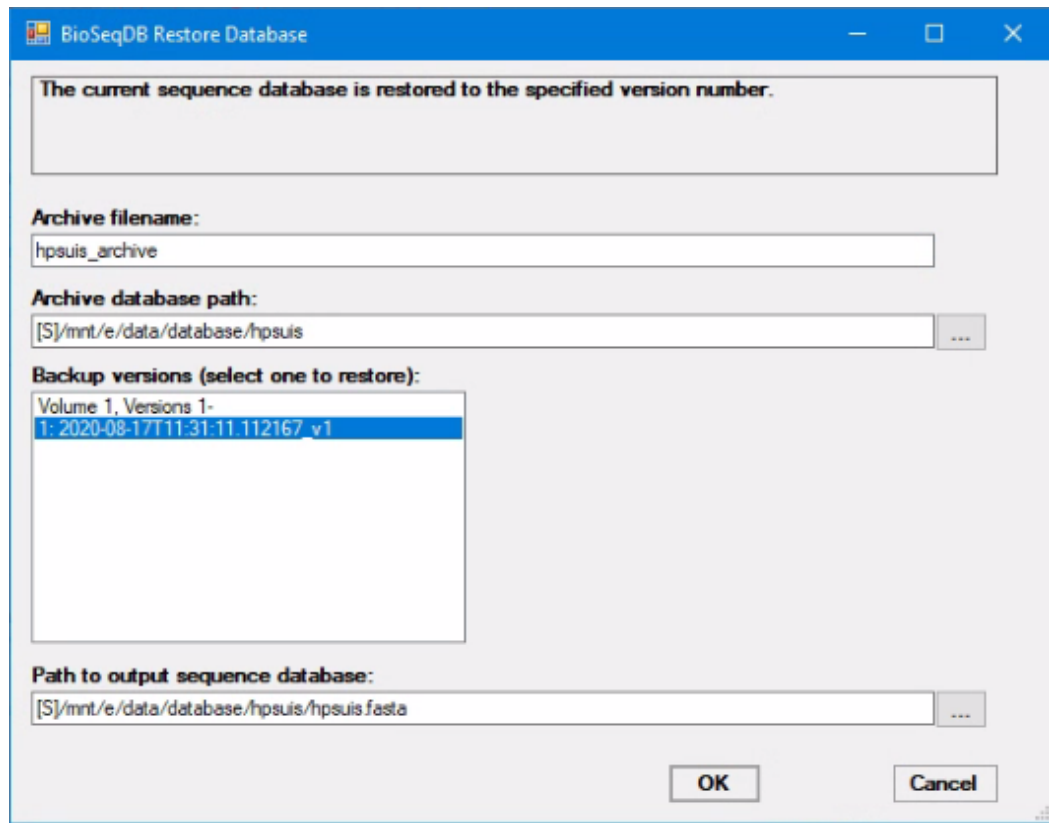
A sequence database may be deleted by clicking on the red 'X' to the right of the selected sequence database name. A confirmation prompt appears to verify that the intention is to delete the database.

6. Backup database



BioSeqDB uses a software tool called Kipper to create an incremental backup of the currently selected database. Each time the Backup function is invoked, only the incremental *changes* to the database are recorded, and a new date/time stamped version of the database is created. The versions are listed on the left of the main dialog. This approach results in a very space-efficient means of creating multiple versions of the database. It is recommended that a backup be performed whenever a significant number of changes to the database have taken place.

7. Restore database

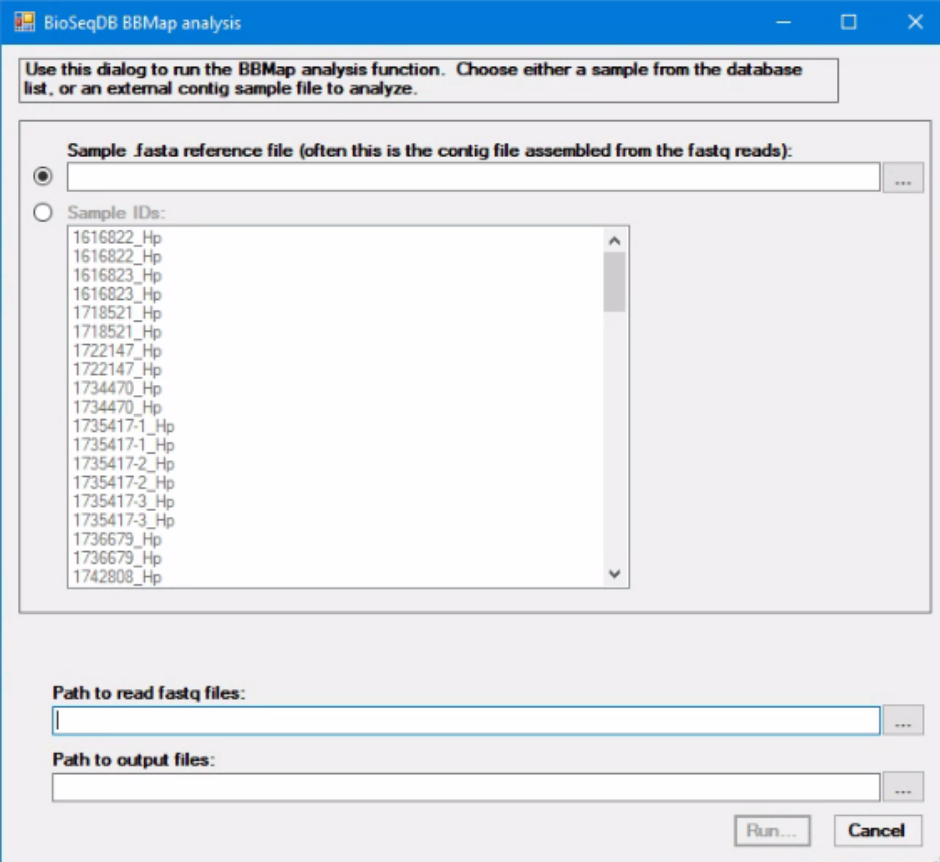


The Restore function can restore the currently selected database to any point in time from the list of backup database versions available. Select the version to restore by clicking on the desired version in the 'Backup versions' list. The currently selected database is replaced by the restored database. If you are restoring an older version of the database but want to later return to the currently selected database, be sure to take a backup first so that later you can restore the current version.

Analysis Functions

Various analysis functions are available which run either in conjunction with data from the sequence databases or with external data. The functions all run as background tasks when activated. Each of the analysis dialogs has an optional memo field to help identify the function when it runs in the background.

0. BBMap

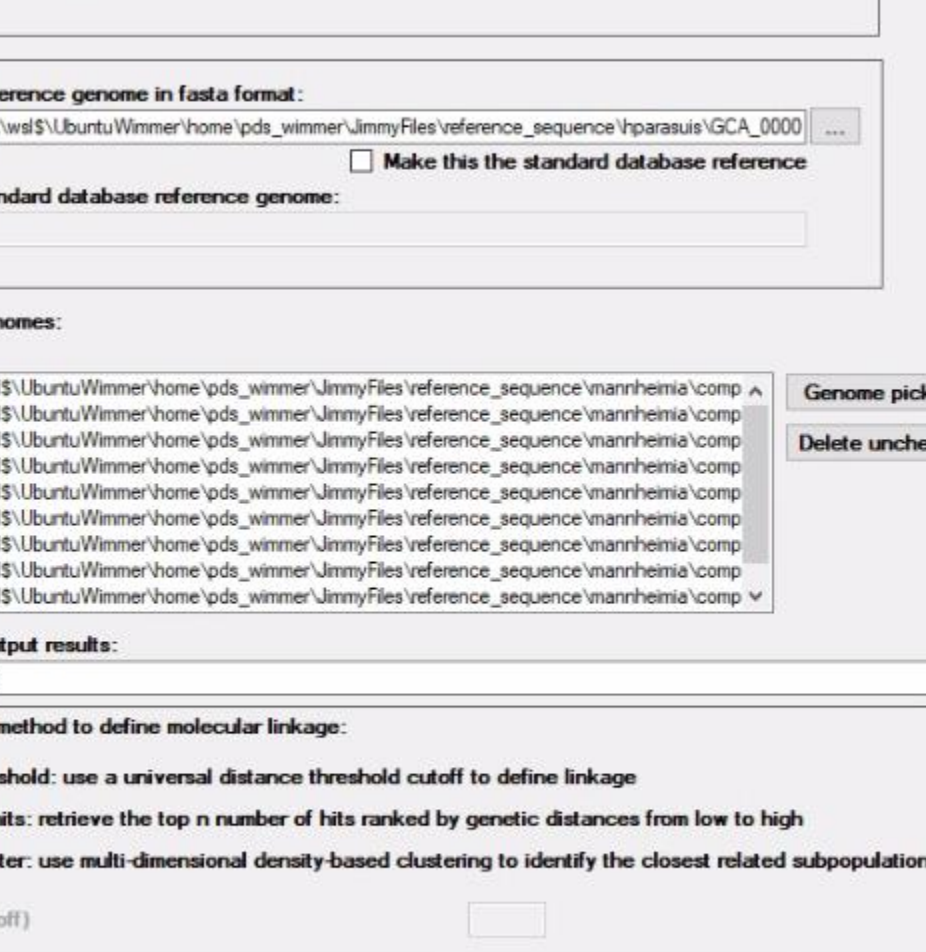


The image shows a Windows-style dialog box titled "BioSeqDB BBMap analysis". At the top, a text box contains the instruction: "Use this dialog to run the BBMap analysis function. Choose either a sample from the database list, or an external contig sample file to analyze." Below this, there are two radio buttons. The first is selected and is labeled "Sample .fasta reference file (often this is the contig file assembled from the fastq reads):". It is followed by an empty text field with a browse button "...". The second radio button is labeled "Sample IDs:" and is followed by a list box containing the following sample IDs: 1616822_Hp, 1616822_Hp, 1616823_Hp, 1616823_Hp, 1718521_Hp, 1718521_Hp, 1722147_Hp, 1722147_Hp, 1734470_Hp, 1734470_Hp, 1735417-1_Hp, 1735417-1_Hp, 1735417-2_Hp, 1735417-2_Hp, 1735417-3_Hp, 1735417-3_Hp, 1736679_Hp, 1736679_Hp, and 1742808_Hp. Below the list box, there are two more text fields. The first is labeled "Path to read fastq files:" and the second is labeled "Path to output files:". Both have empty text fields and browse buttons "...". At the bottom right, there are two buttons: "Run..." and "Cancel".

The BBMap analysis takes as input a reference .fasta file from either an external contig file or a sample from within the selected database. Also specify a path to the .fastq files to match against the reference, as well as a path to where the output results are to be stored.

1. Build tree

The Build tree function is one of the key phylogenetic analysis tools available. The purpose of this function is not limited to just the phylogenetic reconstruction of a set of query genomes, but it also searches the query genomes against database sequences to identify historical strains that demonstrate molecular linkage to query genomes. This functionality enables rapid identification of sequences which likely share lineage origins with query sequences, enabling one to infer epidemiological origins and unknown phenotypic characteristics. To run the phylogenetic analysis, you must specify a reference genome for SNP typing, a set of query genomes, an output path for the results, and a method to define molecular linkage/close relatedness. More information below.



BioSeqDB Build Tree Analysis

Construct a cgSNP tree from closest genomic neighbours in the database.

Reference genome in fasta format:

☒ [S]\wsl\$\Ubuntu\Wimmer\home\pds_wimmer\JimmyFiles\reference_sequence\hparasuis\GCA_0000 ...

☐ Make this the standard database reference

Standard database reference genome:

Query genomes:

☐

☒ [S]\wsl\$\Ubuntu\Wimmer\home\pds_wimmer\JimmyFiles\reference_sequence\mannheimia\comp ^

☒ [S]\wsl\$\Ubuntu\Wimmer\home\pds_wimmer\JimmyFiles\reference_sequence\mannheimia\comp

☒ [S]\wsl\$\Ubuntu\Wimmer\home\pds_wimmer\JimmyFiles\reference_sequence\mannheimia\comp

☒ [S]\wsl\$\Ubuntu\Wimmer\home\pds_wimmer\JimmyFiles\reference_sequence\mannheimia\comp

☒ [S]\wsl\$\Ubuntu\Wimmer\home\pds_wimmer\JimmyFiles\reference_sequence\mannheimia\comp

☒ [S]\wsl\$\Ubuntu\Wimmer\home\pds_wimmer\JimmyFiles\reference_sequence\mannheimia\comp

☒ [S]\wsl\$\Ubuntu\Wimmer\home\pds_wimmer\JimmyFiles\reference_sequence\mannheimia\comp

☒ [S]\wsl\$\Ubuntu\Wimmer\home\pds_wimmer\JimmyFiles\reference_sequence\mannheimia\comp

☒ [S]\wsl\$\Ubuntu\Wimmer\home\pds_wimmer\JimmyFiles\reference_sequence\mannheimia\comp

Genome picker...

Delete unchecked

Path to output results:

[S]\C:\Temp ...

Select method to define molecular linkage:

☐ threshold: use a universal distance threshold cutoff to define linkage

☐ tophits: retrieve the top n number of hits ranked by genetic distances from low to high

☒ cluster: use multi-dimensional density-based clustering to identify the closest related subpopulation

(no cutoff)

☒ Use FastTree analysis (faster) Number of threads:

Memo (to help identify task in notifications):

OK Cancel

Either of two reference genomes can be used, one that is a standard reference associated with the currently selected database, and the other chosen from anywhere in the file system.

Reference genome in fasta format:

☐ C:\data\mannheimia\mannheimia_USDA-ARS-USMARC-186.fasta

☐ Make this the standard database reference

Standard database reference genome:

☒ C:\data\database\mhl\mannheimia_USDA-ARS-USMARC-186.fasta

Use the radio buttons to select which reference genome will be used.

The only way to record a standard reference genome is to select it from the file system and check the checkbox that says, 'Make this the standard database reference'. When the Build tree function runs, the reference genome from the file system is copied to the currently selected database path as the standard reference genome.

Use the 'Genome picker' to select the series of query genomes. All checked genomes are included in the analysis. To remove a genome from the list, uncheck it and select 'Delete unchecked'. This removes ALL genomes that are unchecked in the list.

To identify close genomic neighbors in the database, a method to define molecular linkage must be selected. There are three different methods available:

1. **threshold** – set a universal cutoff value in which sequences below the given distance threshold are defined as close genomic neighbors
2. **tophits** – the top n number of database sequences ranked by genetic distances from low to high are defined as close genomic neighbors
3. **cluster** – a non-parametric method that uses density-based clustering to dynamically characterize subpopulations based on the genetic distances between the query and database sequences. The subpopulation with the shortest distance to the query is selected as close genomic neighbors.

There are other optional parameters that may be specified. A checkbox to specify 'Use Fast Tree analysis' results in an analysis that runs much faster but is slightly less accurate. Another option is to adjust the distance threshold values for the linkage methods. This only applies to 'threshold' and 'tophits' linkage methods, as the 'cluster' method is non-parametric.

The output of the Build tree function is two files, tree.nwk and metadata_microreact.csv. The tree.nwk file consists of the tree description in Newick format, and the metadata_microreact.csv file that consists of a list of sample IDs mapped to the status attribute. The status attribute indicates whether a given sample in the tree is a query or a local database sequence. You can append more sample attributes to this file such as AMR profiles or serotypes to highlight certain

trends in the phylogentic tree. Note that the metadata_microreact.csv is only designed to annotate trees visualized on the Microreact Internet platform. Upon analysis completion, the Build tree function automatically invokes the Dendroscope program to display the tree graphically, but the files may also be manually uploaded to the <https://microreact.org/upload> website for an alternative method of displaying the tree.

Use the Memo field to enter identifying information to facilitate tracking the background task in the Notifications dialog.

2. Influenza A pipeline

Construct the consensus sequence of the full-length Influenza A genome from Nanopore reads. The pipeline uses Centrifuge to determine the segment identity (segment 1-8) of each read from a single fastq file and subsequently performs multiple rounds of genome polishing steps per segment using medaka and racon.

Fastq files representing Influenza A samples:

☐ (all/none)

Sample ID	Path
<input type="checkbox"/> PDS2100454-...	[S]E:\data\Dreamer\InfA_Sequencing\InfA_5samples_Flowcell_18Jan21DP\InfA_5e
<input type="checkbox"/> PDS2100454-...	[S]E:\data\Dreamer\InfA_Sequencing\InfA_5samples_Flowcell_18Jan21DP\InfA_5e
<input type="checkbox"/> PDS2100454-...	[S]E:\data\Dreamer\InfA_Sequencing\InfA_5samples_Flowcell_18Jan21DP\InfA_5e
<input type="checkbox"/> PDS2100454-...	[S]E:\data\Dreamer\InfA_Sequencing\InfA_5samples_Flowcell_18Jan21DP\InfA_5e
<input type="checkbox"/> PDS2100454-...	[S]E:\data\Dreamer\InfA_Sequencing\InfA_5samples_Flowcell_18Jan21DP\InfA_5e
<input type="checkbox"/> PDS2100454-1	[S]E:\data\Dreamer\InfA_Sequencing\InfASeq_Flongle_5samples_8Jan20_DP\PD
<input type="checkbox"/> PDS2100454-2	[S]E:\data\Dreamer\InfA_Sequencing\InfASeq_Flongle_5samples_8Jan20_DP\PD
<input type="checkbox"/> PDS2100454-3	[S]E:\data\Dreamer\InfA_Sequencing\InfASeq_Flongle_5samples_8Jan20_DP\PD
<input type="checkbox"/> PDS2100454-4	[S]E:\data\Dreamer\InfA_Sequencing\InfASeq_Flongle_5samples_8Jan20_DP\PD
<input checked="" type="checkbox"/> PDS2100454-5	[S]E:\data\Dreamer\InfA_Sequencing\InfASeq_Flongle_5samples_8Jan20_DP\PD

Fastq file picker...
Delete unchecked

Path to Centrifuge database for taxonomic and segment classification:
[S]E:\data\database\centrifuge\viralRefSeq_InfA_custom

Centrifuge database name: viralRefSeq_InfA_custom

Path to output result files, including final consensus sequence in subfolder <SampleID>\consensus\
[L]C:\\Temp

☐ Adaptor trimming by porechop Number of threads: 12 Flowcell chemistry model: r9

Segments to assemble: ☒ All ☐ 1 ☒ 2 ☒ 3 ☐ 4 ☐ 5 ☐ 6 ☐ 7 ☐ 8

Memo (to help identify task in notifications):

Run Cancel

The purpose of this pipeline is to generate the consensus sequence of all eight genomic segments of Influenza A from raw Nanopore sequences that were amplified using universal PCR primers. As such, the required inputs are .fastq files which can be selected using the “Fastq file picker”. Given that the raw Nanopore reads of each sample are partitioned across multiple fastq files, we must select the **DIRECTORY** that contains the individual fastq files of each sample rather than selecting individual fastq files. A sample ID must be given to keep track of each sample and the name of the selected directory is used as the sample ID as default but can be changed if necessary.

This pipeline uses Centrifuge to perform taxonomic classification and alignment of each read to the genomic segments of Influenza A. The second step is to indicate the **path to the DIRECTORY** that contains the Centrifuge database for taxonomic classification. There is an analysis-ready Centrifuge database in E:\data\database\centrifuge\viralRefSeq_InfA_custom.

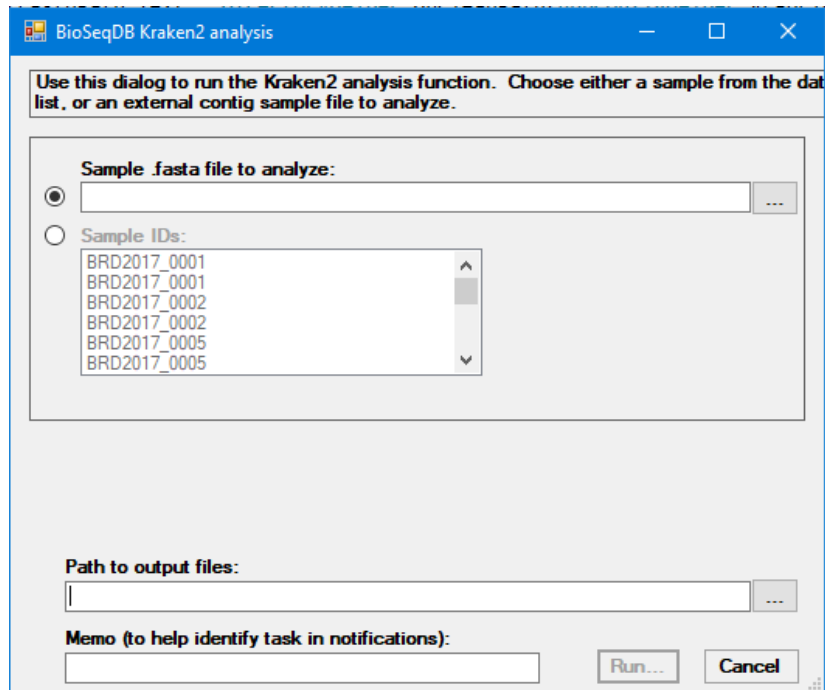
The third required input is the path to the output directory where the consensus sequences of each sample will be stored. The output of each sample will be organized into individual subdirectories in the specified output directory. If the output path is on the local computer, the results are copied to that folder on the local machine.

Raw read adaptor and barcode trimming by porechop can be enabled as an optional step; however, it is strongly recommended to avoid having adaptor sequences embedded in the final consensus sequence. The number of threads can also be adjusted to request a higher number of cores to run the analysis.

Influenza A virus has eight segments. All eight segments can be assembled when the analysis is run, or any subset of segments may be selected if there is a need to focus on only one segment or a few segments. The flowcell chemistry model, R9 or R10 can also be specified.

Use the Memo field to enter identifying information to facilitate tracking the background task in the Notifications dialog.

3. Kraken2



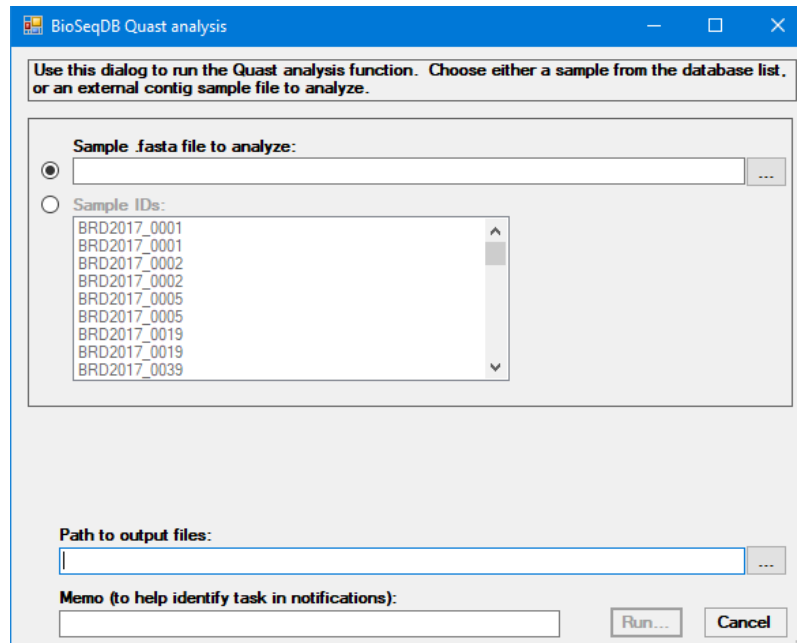
The dialog box is titled "BioSeqDB Kraken2 analysis". It contains the following elements:

- Instructions:** "Use this dialog to run the Kraken2 analysis function. Choose either a sample from the data list, or an external contig sample file to analyze."
- Sample .fasta file to analyze:** A radio button is selected next to a text field with a browse button "...".
- Sample IDs:** A radio button is unselected next to a list box containing the following IDs:
 - BRD2017_0001
 - BRD2017_0001
 - BRD2017_0002
 - BRD2017_0002
 - BRD2017_0005
 - BRD2017_0005
- Path to output files:** A text field with a browse button "...".
- Memo (to help identify task in notifications):** A text field.
- Buttons:** "Run..." and "Cancel".

The Kraken2 analysis takes as input a reference .fasta file from either an external contig file or a sample from within the selected database. Also specify a path to where the Kraken2 output results are to be stored.

Use the Memo field to enter identifying information to facilitate tracking the background task in the Notifications dialog.

4. Quast



The Quast quality assessment analysis takes as input a reference .fasta file from either an external contig file or a sample from within the selected database. Also specify a path to where the Quast output results are to be stored.

Use the Memo field to enter identifying information to facilitate tracking the background task in the Notifications dialog.

5. Salmonella Serotyping

The screenshot shows a Windows application window titled "BioSeqDB Salmonella Serotyping". At the top, a text box explains the pipeline: "The snakemake pipeline functions to predict Salmonella serotypes from raw Nanopore sequencing data. The pipeline uses flye to assemble the long reads, medaka for genome polishing, and SISTR for in-silico serotyping." Below this, the "Fastq files representing Salmonella samples:" section contains a checkbox for "(all/none)" and a table with two columns: "Sample ID" and "Path". The table has one row with "fastq" checked in the "Sample ID" column and "[S]D:\data\BRD_24_26Mar2020\fastq\" in the "Path" column. To the right of the table are buttons for "Fastq file picker..." and "Delete unchecked". Below the table is a "Path to output folder" section with a text box containing "[L]C:\\Temp" and a browse button "...". Further down is a checkbox for "Adaptor trimming by porechop". At the bottom, there is a "Memo (to help identify task in notifications):" section with a text box, and "Run" and "Cancel" buttons.

Sample ID	Path
<input checked="" type="checkbox"/> fastq	[S]D:\data\BRD_24_26Mar2020\fastq\

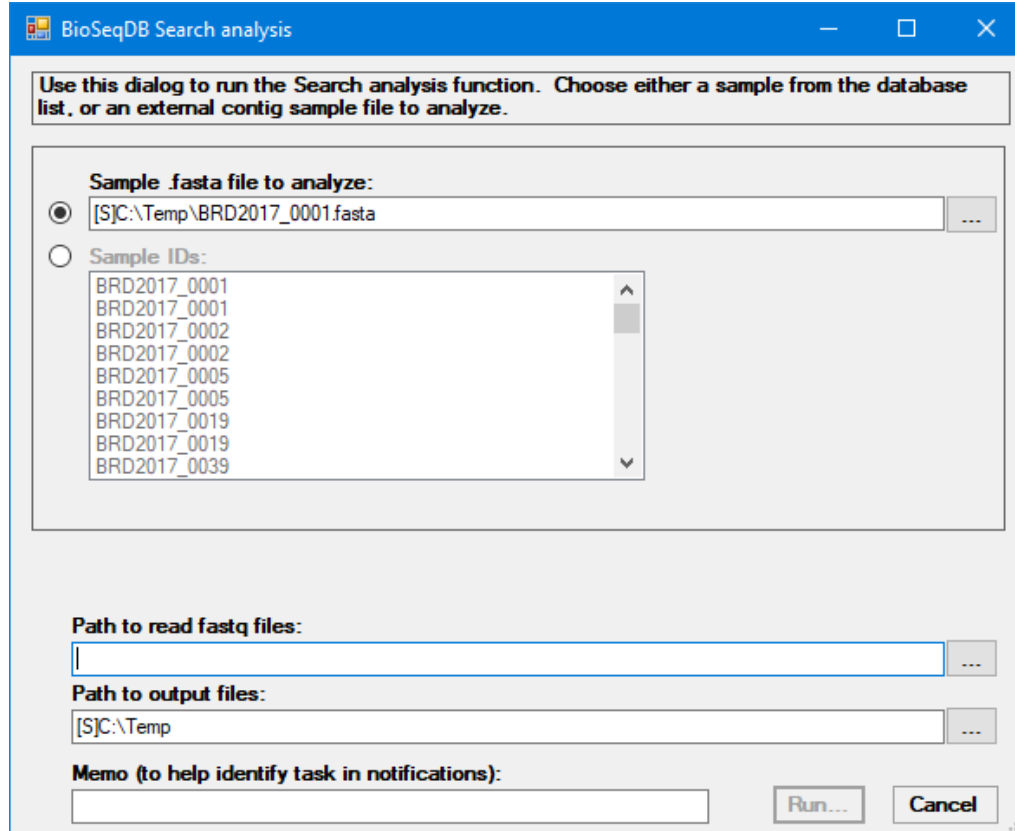
The purpose of this pipeline is to perform in-silico Salmonella serotyping from raw Nanopore sequences. As such, the required inputs are fastq files which can be selected using the “Fastq file picker”. Given that the raw Nanopore reads of each sample are partitioned across multiple fastq files, we must select the **DIRECTORY** that contains the individual fastq files of each sample rather than selecting individual fastq files. A sample ID must be given to keep track of each sample and the name of the selected directory is used as the sample ID as default but can be changed if necessary.

The second required input is the path to the output directory where the serotyping results will be written to. The serotyping results of all the samples selected in a single analysis run are aggregated into a single output file named “sistr_res_aggregate{time}.csv”. If the output path is on the local computer, the results are copied to that folder on the local machine and automatically opened.

Raw read adaptor and barcode trimming by porechop can be enabled as an optional step. The number of threads can also be adjusted to request a higher number of cores to run the analysis.

Use the Memo field to enter identifying information to facilitate tracking the background task in the Notifications dialog.

6. Search



The image shows a Windows-style dialog box titled "BioSeqDB Search analysis". It contains a message box at the top stating: "Use this dialog to run the Search analysis function. Choose either a sample from the database list, or an external contig sample file to analyze." Below this, there are two main options for selecting a sample to analyze. The first option, "Sample .fasta file to analyze:", is selected with a radio button and shows a text field containing "[S]C:\Temp\BRD2017_0001.fasta" with a browse button "...". The second option, "Sample IDs:", is unselected and shows a list box containing the following IDs: BRD2017_0001, BRD2017_0001, BRD2017_0002, BRD2017_0002, BRD2017_0005, BRD2017_0005, BRD2017_0019, BRD2017_0019, and BRD2017_0039. Below the list box, there are three more text fields: "Path to read fastq files:" (empty), "Path to output files:" (containing "[S]C:\Temp"), and "Memo (to help identify task in notifications):" (empty). At the bottom right, there are "Run..." and "Cancel" buttons.

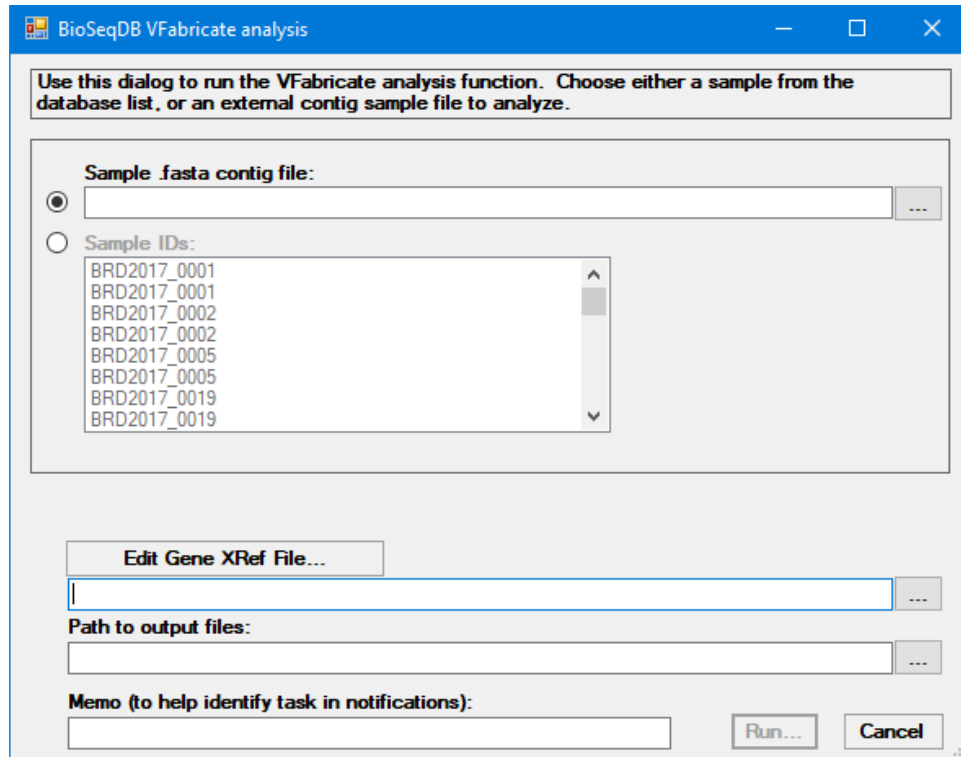
The Search function searches the currently selected database for sequences like a query sequence. You are prompted for a new sample ID that is used to name the output result file, a path to a query sequence in .fasta format, and a path where the output results are stored. If there are no errors, the results are stored in a .txt file and you are prompted to open the file. The result file contains a list of up to 50 sequences that are like the query sequence in descending order of similarity.

Two options are available to qualify the search. The results may be filtered by a maximum distance cutoff, ranging from 0 to 1 with a default of 1. The number of threads may also be specified to speed up the search.

Use the Memo field to enter identifying information to facilitate tracking the background task in the Notifications dialog.

7. VFabricate

The VFabricate analysis uses Abricate to report on the frequency of virus genes in each sample. The sample comes from either an external contig file or a sample from within the selected database. Also specify a path to where the VFabricate output results are to be stored.



The dialog box is titled "BioSeqDB VFabricate analysis". It contains a text box with the instruction: "Use this dialog to run the VFabricate analysis function. Choose either a sample from the database list, or an external contig sample file to analyze." Below this, there are two radio buttons. The first is labeled "Sample .fasta contig file:" and is selected. The second is labeled "Sample IDs:". Below the second radio button is a list box containing the following sample IDs: BRD2017_0001, BRD2017_0001, BRD2017_0002, BRD2017_0002, BRD2017_0005, BRD2017_0005, BRD2017_0019, and BRD2017_0019. Below the list box is a button labeled "Edit Gene XRef File...". Below this button is a text box for the "Path to output files:". Below the text box is a button labeled "Run..." and a button labeled "Cancel". At the bottom, there is a text box labeled "Memo (to help identify task in notifications):".

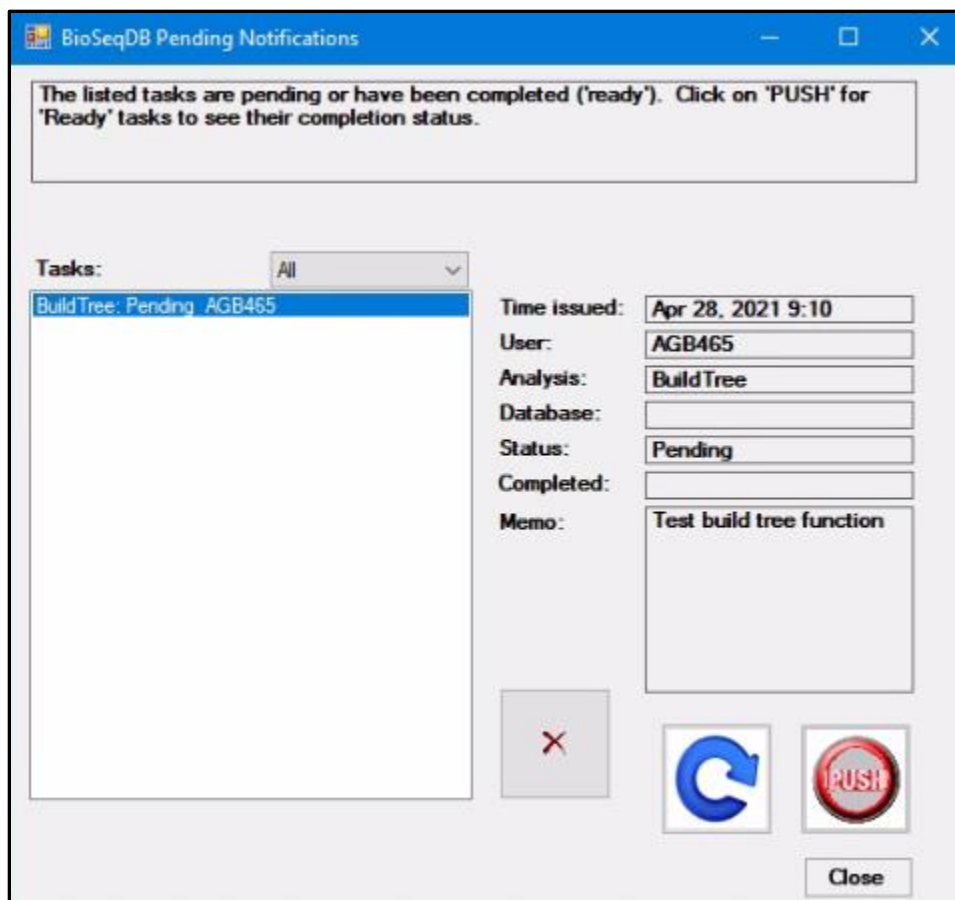
The Gene XRef file describes the target gene name prefixes that are identified in the Abricate results from each of the online virus databases (Card, Ecoli_VF, ECOH, VFDB, NCBI, ResFinder, PlasmidFinder, Argannot). Each line contains a gene name prefix and a description, separated by a comma, as follows:

```
apx, Actinobacillus pleuropneumoniae toxin
apf, Type 4 fimbriae
cpx, Capsular polysaccharide
pap, P fimbriae
sfa, S fimbriae
cnf1, Cytotoxic necrotizing factor 1
hly, Hemolysin
iuc, Aerobactin
iut, Aerobactin
```

Use the Memo field to enter identifying information to facilitate tracking the background task in the Notifications dialog.

Running in the background

The Assemble function and the analysis functions are tasks with the potential to take a long time to complete. Rather than cause BioSeqDB to wait until they complete, they are scheduled to run in the background while other BioSeqDB functions can be performed. The tasks running in the background have a status of 'Pending'. When they complete, their status changes to 'Ready'. When a task becomes 'Ready', the user is notified through the 'Pending Notifications' dialog.



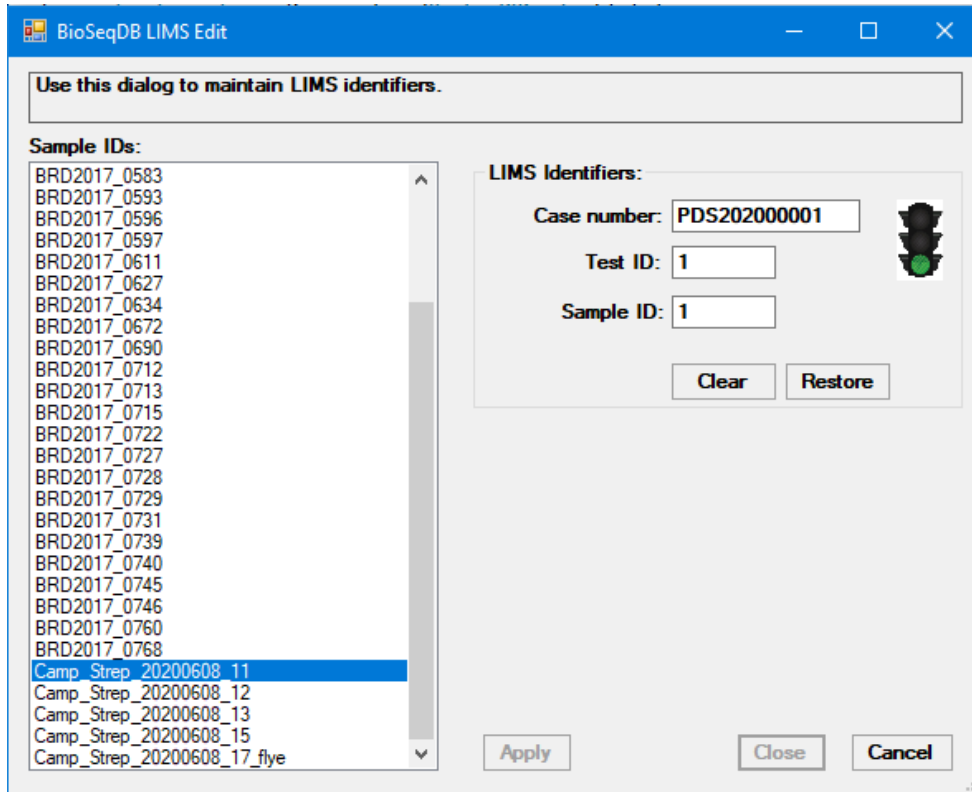
The details of the selected task in the Tasks list are displayed on the right of the dialog. When a task completes, the Status changes to 'Ready', and the background of the Status field changes to orange. The time of completion is also reported.

A task, whether the status is 'Pending' or 'Ready', may be deleted from the list at any time. You can also refresh the list. When a task status changes to 'Ready', the 'Push' button becomes enabled. By clicking on the 'Push' button, the result of the task is displayed, and the task is removed from the task list. For example, for the Build tree function, if successful, the result of the function is displayed as a phylogenetic tree in Dendroscope. For the Assemble function, the results of the assembly are stored in the E:\data\staging\<sample> folder.

The main dialog has a 'Pending Notifications' command button that opens the Pending Notifications dialog when clicked. The text of the command button also displays the number of Pending and Ready tasks in the task list of the Pending Notifications dialog.

Editing LIMS identifiers

Any of the existing sequence samples in the currently selected sequence database may be assigned to the LIMS samples via the BioSeqDB LIMS edit dialog:



The dialog box is titled "BioSeqDB LIMS Edit". It contains a message at the top: "Use this dialog to maintain LIMS identifiers." Below this, there are two main sections. The left section, labeled "Sample IDs:", contains a list of sample IDs. The right section, labeled "LIMS Identifiers:", contains three input fields: "Case number:", "Test ID:", and "Sample ID:". Below these fields are two buttons: "Clear" and "Restore". At the bottom of the dialog are three buttons: "Apply", "Close", and "Cancel".

Sample IDs:

- BRD2017_0583
- BRD2017_0593
- BRD2017_0596
- BRD2017_0597
- BRD2017_0611
- BRD2017_0627
- BRD2017_0634
- BRD2017_0672
- BRD2017_0690
- BRD2017_0712
- BRD2017_0713
- BRD2017_0715
- BRD2017_0722
- BRD2017_0727
- BRD2017_0728
- BRD2017_0729
- BRD2017_0731
- BRD2017_0739
- BRD2017_0740
- BRD2017_0745
- BRD2017_0746
- BRD2017_0760
- BRD2017_0768
- Camp_Strep_20200608_11**
- Camp_Strep_20200608_12
- Camp_Strep_20200608_13
- Camp_Strep_20200608_15
- Camp_Strep_20200608_17_flye

LIMS Identifiers:

Case number: PDS202000001

Test ID: 1

Sample ID: 1

Clear Restore

Apply Close Cancel

Select a sequence sample from the list on the left, and use the edit are on the right to add, modify or clear any associated values. The traffic light icon indicates when value combinations represent valid values. When 'Apply' is enabled, this indicates that there are potential changes in values that you are editing. When you click on 'Apply', those changes are recorded permanently. The 'Apply' button is only enabled when the LIMS identifier values are valid.

For BioSeqDB administrators

- Control file appsettings.json

BioSeqDB has its own repository of variables and values that it uses to keep track of global definitions, preferences and selections, currently selected databases, and outstanding tasks. The file name is appsettings.json and it resides in the same folder as the BioSeqDB service executable on the server (C:\BioSeqDB\Service\appsettings.json). The data structure is mapped by the BioSeqDBConfig class in the BioSeqDBConfig.cs source file in BioSeqDB.

There is also an equivalent appsettings.json file for each individual user to remember individual user preferences and options. The naming convention for these appsettings files is appsettings_<username>.json.

The basic structure of the appsettings.json file is:

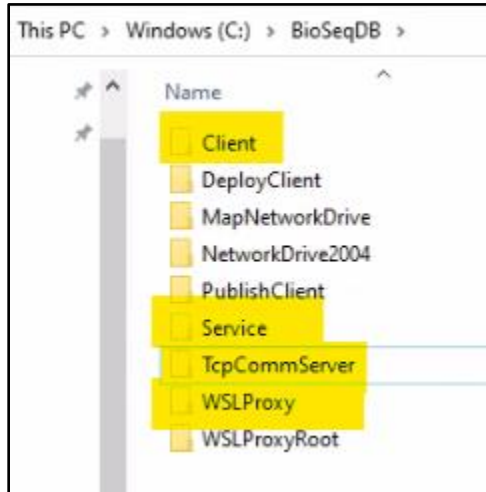
```
{
  <global settings>
  "seqDBs":
  {
    <list of DBs and their properties>
  }
  "Users":
  {
    <list of users and their properties>
  }
  "Tasks":
  {
    <list of tasks and their properties (empty if no outstanding tasks)>
  }
}
```

Normally you should not make any changes to the values in the appsettings.json file. However, sometimes for diagnostic purposes it is useful to examine the values recorded in this file.

If updating the executable files for BioSeqDB service, take care not to overwrite the appsettings.json file or the appsettings files for any of the users.

- Updating the executables

There are several components to the BioSeqDB platform in addition to the BioSeqDB client.



Each of these components has a folder under the C:\BioSeqDB folder on WIMMER from which they run. These are the Client, Service, TcpCommServer and WSLProxy subfolders. Note that this is also the location of the WSLProxyRoot folder, under which are stored all the temporary data related to individual users. These data may be deleted if large amounts are accumulating.

The BioSeqDB service runs under the Network Authority account and is responsible for performing all BioSeqDB service requests, including translating any Linux command requests for WSLProxy to

perform. The service uses IP address 8390, which must be opened in the firewall.

Normally when updating the service executable, stop the service from the Windows Services function, replace the changed .exe, .dll and .pdb files, then start the service. These files are found in the BioSeqDBSolution\BioSeqService\bin\Debug folder of the development environment.

Do not make any changes to any config or appsettings files. Note that diagnostic information is logged in the 'logs' subfolder of the Service subfolder. This is generally true for all components.

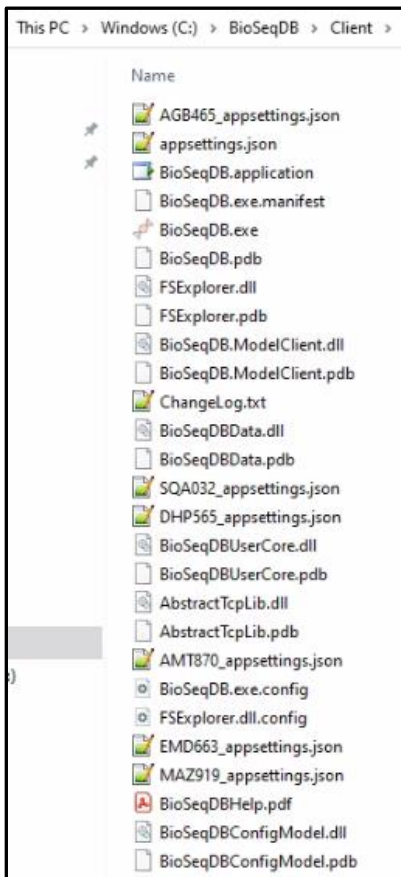
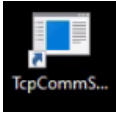
The WSLProxy service is started automatically on WIMMER and responds to any Linux requests passed to it from the BioSeqDB service. It runs under the local Administrator account, and so can perform services that the BioSeqDB service is unable to perform, including launching Linux commands. Like the BioSeqDB service, stop the WSLProxy service by closing the WSLProxy dialog, replace the changed .exe, .dll and .pdb files, then restart the service by clicking on the WSLProxy icon.



The TcpCommServer is a specialized file transfer service that uses the TCP protocol to transfer files between the local computer and the server (WIMMER). It starts automatically on WIMMER and communicates with the Explorer in BioSeqDB whenever file transfers are requested. TcpCommServer is capable of high-speed, parallel data transfer operations, making it a highly scalable file transfer service.

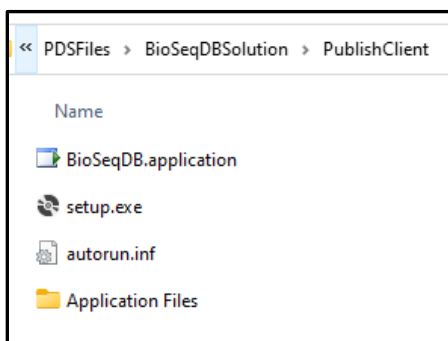
TcpCommServer uses port 22490 for file transfers, which must be open in the firewall.

Like the WSLProxy service, stop the TcpCommServer by closing its dialog on WIMMER, replace the changed .exe, .dll and .pdb files in the TcpCommServer folder, and restart the service by clicking on the TcpCommServer icon.



The client component has two methods of update, depending on whether the client is running from the WIMMER server or if it is being deployed to a client computer via ClickOnce.

On WIMMER, after making sure the BioSeqDB client is shut down, replace the BioSeqDB.application, BioSeqDB.exe.manifest, .exe, .dll, and .pdb files. If necessary, also replace the ChangeLog.txt and BioSeqDBHelp.pdf files. Do not touch the .json files or any config files.



The second method for deploying the client involves doing a 'Publish' of BioSeqDB in the development environment, then copying the latest version of the published folder in the 'Application Files' folder to the W:\PDS\Apps\BioSeqDB\Application Files folder on the jade.usask.ca share. This can be accessed on the W: drive on the PDS2014006 jump box computer.

Once that published subfolder is in place, all that is required is to replace the BioSeqDB.application file from the C:\PDSFiles\BioSeqDBSolution\PublishClient folder on the development machine to the W:\PDS\Apps\BioSeqDB folder. This ensures that the manifest file is in sync with the application folder of the current update.

Once this is in place, users can initiate the execution process by running the setup.exe in the same folder, assuming they can access the W: drive from their own machine. There should be a permanent mapping to the W: drive on the client computer because when the client is started it automatically checks for and downloads an update of the application if one is available.

Note that this same ClickOnce deployment mechanism is available for running remotely via VPN. However, in that case there is no mapping to the W: drive, so the whole W:\PDS\Apps\BioSeqDB folder must be physically transported to the remote computer for deployment. Automatic updates are not available in this scenario, but this is a small price to pay for this capability.

- User ID management

Currently BioSeqDB uses a quite simple user management system. There is a list of usernames and passwords in the appsettings.json file that represents valid users of BioSeqDB. These names are maintained manually. This is one exception where the appsettings.json file must be edited manually. Simply make whatever changes are necessary to identify the BioSeqDB users and save the changes. The passwords must be manually encrypted for the initial setup, but the user may change their own passwords once they are registered. There is an encryption utility available in the development sandbox project.

Take care to conform to the syntax of Json to avoid errors at startup.

- Source control

The BioSeqDB source is stored on GitHub at <https://github.com/ArnieBerg/BioSeqDBSolution>. The source for the Linux seqdb scripts is stored at <https://github.com/jimmyliu1326/seqdb>.

- Updating this Help file

This file is called BioSeqDBHelp.docx in the \PDSFiles\BioSeqDBSolution\BioSeqDB folder in the development environment. After editing it to reflect any recent changes, save it as a PDF file with the same name in the same folder. The PDF file will automatically be included in the next ClickOnce deployment.

FAQ

- 1. If I have a standard reference genome defined for the currently selected database, and I do a backup, is the reference genome backed up as well?**

No, at this point the reference genome is not backed up.

- 2. The 'BioSeqDB LIMS edit' dialog allowed me to enter a case number/test ID/sample ID that does not exist in the LIMS. Is that correct?**

Yes, at this point there is no attempt to cross-reference the values you enter with actual LIMS identifiers.

- 3. Which molecular linkage method should I choose for Build Tree?**

This really depends on how specific you want to define molecular linkage. If you are looking to identify epidemiological linked sequences, you may want to use a strict distance cutoff such that you have 100% confidence that the sequences you identify are very closely related to the queries. However, you must keep in mind that there might be sequences which happen to miss the distance threshold by a small margin and could potentially be highly relevant.

If you are unaware of a suitable distance threshold to define close relatedness for your query genomes, then using `tophits` or `cluster` methods are recommended. However, you should also keep in mind that with the `tophits` method, because it is a rank-based method, the identified sequences could in fact be highly distant from your queries or there could be a high abundance of closely related sequences, but some were excluded for not meeting the top-ranking threshold.

Hence, to avoid the exclusion of relevant strains due to user-defined thresholds, we introduced the cluster method that attempts to optimize the parameters that should be used to define subpopulations closely related to the queries. However, with the `cluster` method, the approach relies on the existence of an underlying population structure within the database. Consequently, this approach is not ideal for small sized databases with limited genetic diversity.

- 4. One of my pending tasks stays pending even though I know it completed. How is this possible?**

This can happen if BioSeqDB is restarted while tasks are pending. It loses track of which tasks are pending and ready. The tasks continue to appear in the notification list, but the status may not be correct. Best practice is to leave BioSeqDB running continuously if possible.

- 5. What if the output of one of the background tasks produces a large amount of standard output information?**

Normally when a task completes, there is a small amount of information reported as to the progress, success and/or error associated with the task. If the amount of information exceeds

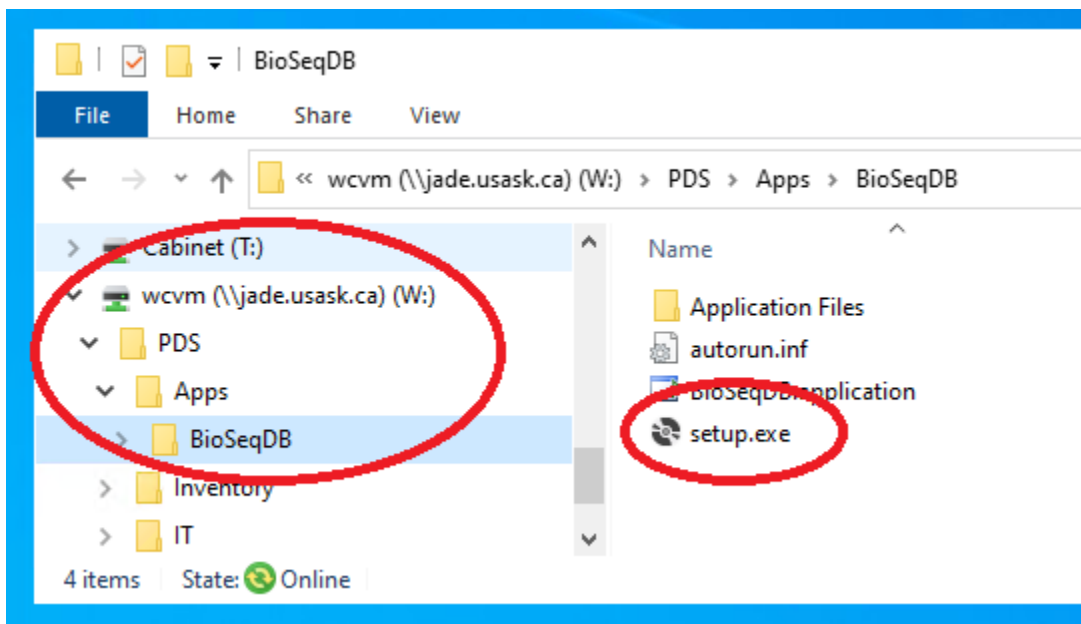
the capacity of the dialog, it is stored in a TaskLog file and opened automatically in the default Windows editor.

Appendix A: Deploying the BioSeqDB client

Before deploying the BioSeqDB client to your computer, ensure the following have been done:

1. The BioSeqDB administrator must register your NSID in BioSeqDB.
2. Make sure C:\Temp exists on your computer.
3. Install Dendroscope and set PathToDendroscope in appsettings for user. If Dendroscope is not installed, the output of the Build Tree function cannot be visualized.
4. Make sure Excel is installed. If Excel is not installed, some output results may not be viewable in Excel format.

The BioSeqDB client deployment is handled by an installation method known as ClickOnce. After installing on your computer for the first time, you will have an entry for BioSeqDB in the Start menu. On subsequent times when you launch BioSeqDB, it will automatically check whether a newer version is available and prompt you to install it.



To access the setup file for BioSeqDB, you should have a mapped network drive to W:\\PDS\\Apps\\BioSeqDB.

Double-click on the **setup.exe** file in the BioSeqDB folder and click on the prompts to start the installation process. Once the installation is complete, the application will start.

Once the application is installed, all you need to do to keep it up to date is to launch it from the start menu. If a newer version is available, you will be prompted to install it.

You may also want to pin the BioSeqDB app to the taskbar. Do this by right-clicking on the BioSeqDB icon on the taskbar while the app is running and selecting 'Pin to taskbar'.

If for some reason you need to uninstall the app, this may be done from 'Programs and Features' in the Control Panel.