

ONTbarcode

A pipeline for MinION based DNA barcoding

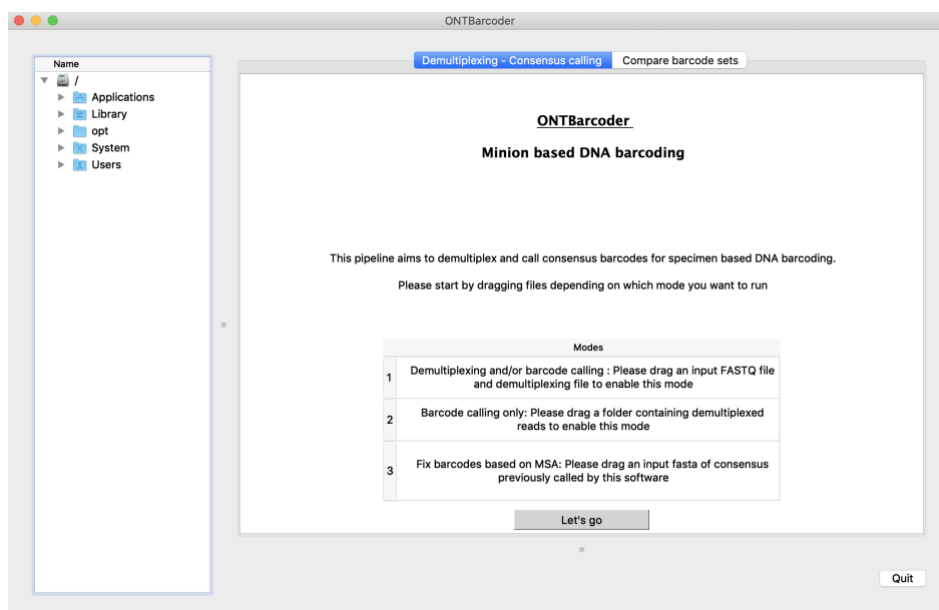
ONTBarcoder (available from <https://github.com/asrivathsan/ONTbarcoder>) is installed by unzipping the folder with the version of the program that supports the operating system on your computer. The folders are available from here:

ONTBarcoder's three modules.

1. Demultiplexing: MinION reads are assigned to specimen-specific bins.
2. Barcode calling: The reads in the specimen-specific bins are used to derive the barcodes based on alignment and consensus calling.
3. Barcode comparison module: Two or more sets of barcodes can be compared.

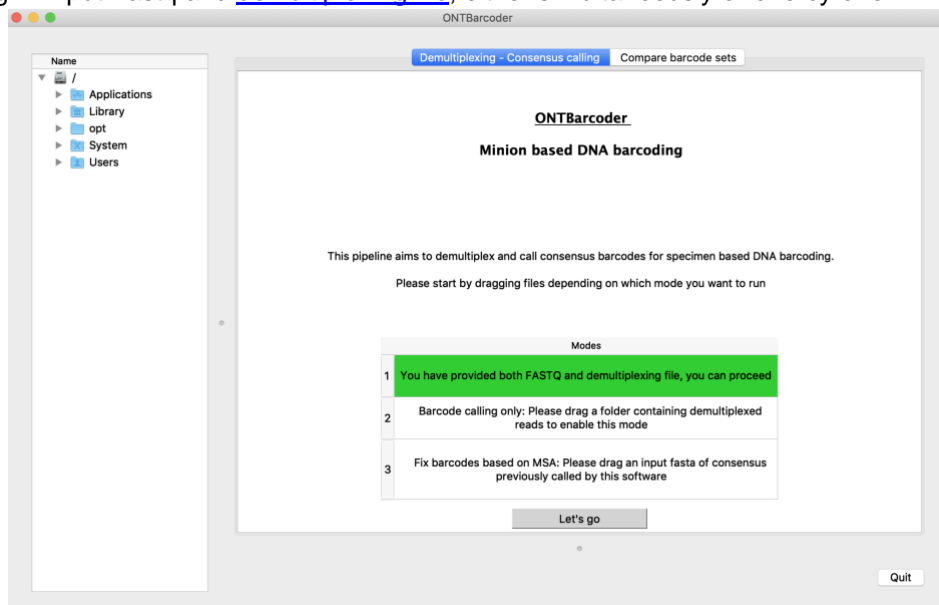
We recommend that you try the pipeline using our DatasetA containing 257 amplicons generated from Flongle. It is available from [here](#). Information to download the various datasets is available [here](#). Prior to running any dataset, ensure you have at least as much space as is required for the input FASTQ file.

1. Open ONTbarcoder



Demultiplexing and barcode calling using the “Demultiplexing-Consensus calling tab”

2. Drag in Input Fastq and [demultiplexing file](#), either simultaneously or one by one

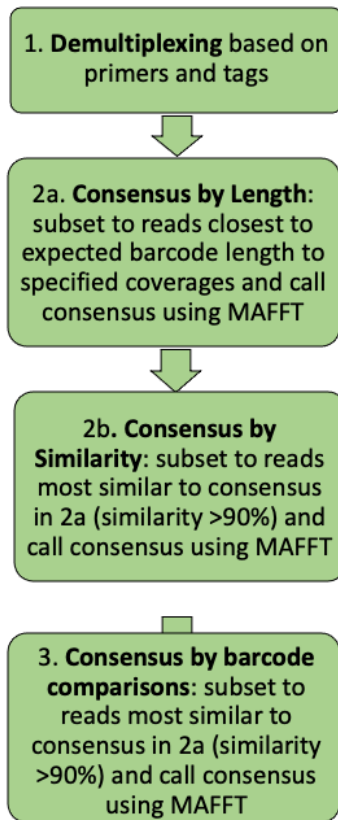


3. Press “Let's go”

4. Configure your run based on the settings described below. DatasetA does not require any changes to the default. The default assumes that you are working on invertebrate barcodes that are 658 bp in length. The remaining settings define which reads are demultiplexed, where the primers are found, how many reads are used in the different steps, and how the consensus barcodes are determined.

BASIC SETTINGS

Analysis workflow



Demultiplexing Settings

Minimum Length: 658

Length of barcode: 658

Window to define product length (length of barcode +/-): 100

Window for primer and tag search: 100

Consensus by length settings

Coverage used: 25,50,100,200,500

Maximum deviation of read length from barcode length: 50

Consensus by similarity settings

Coverage used: 100

General settings for consensus calling

Main consensus calling frequency: 0.3

Range of frequencies to assess: 0.2,0.5

to be examined at step size of 0.05

Genetic Code: 5. The Invertebrate Mitochondrial Code

Steps to run

☒ Consensus by length ☒ Consensus by similarity ☒ Consensus by barcode comparisons

Run

Annotations:

- Reads longer than this are accepted for demultiplexing (points to Minimum Length)
- Expected barcode length, used for QC and determining expected product length (points to Length of barcode)
- Number of reads used for consensus by length. If multiple coverages are specified (delimiter=comma), then iterative mode is used. (points to Coverage used)
- In iterative mode, consensus are called at the lowest coverage. Those barcodes not passing the QC criteria of translation, meeting barcode length and being ambiguity free are passed for barcode calling at higher coverage. (points to Coverage used)
- Number of reads used for consensus by similarity. (points to Coverage used)
- Genetic code (points to Genetic Code)
- Start the run (points to Run button)

ADVANCED SETTINGS

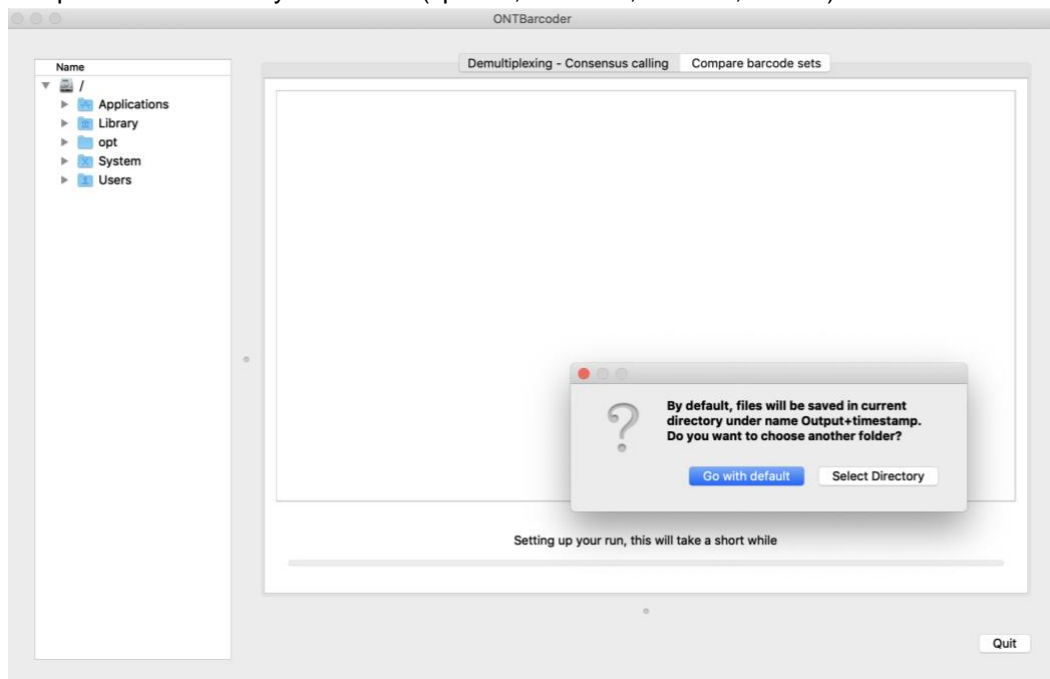
The screenshot shows a software interface for 'ADVANCED SETTINGS' with several sections and input fields. Green arrows point from text boxes to specific settings:

- Demultiplexing Settings**
 - Minimum Length: 658
 - Length of barcode: 658
 - Window to define product length (length of barcode +/-): 100
 - Window for primer and tag search: 100
- Consensus by length settings**
 - Coverage used: 25,50,100,200,500
 - Maximum deviation of read length from barcode length: 50
- Consensus by similarity settings**
 - Coverage used: 100
- General settings for consensus calling**
 - Main consensus calling frequency: 0.3
 - Range of frequencies to assess: 0.2,0.5
 - to be examined at step size of: 0.05
- Genetic Code:** 5. The Invertebrate Mitochondrial Code
- Steps to run:**
 - ☒ Consensus by length
 - ☒ Consensus by similarity
 - ☒ Consensus by barcode comparisons
- Run** button

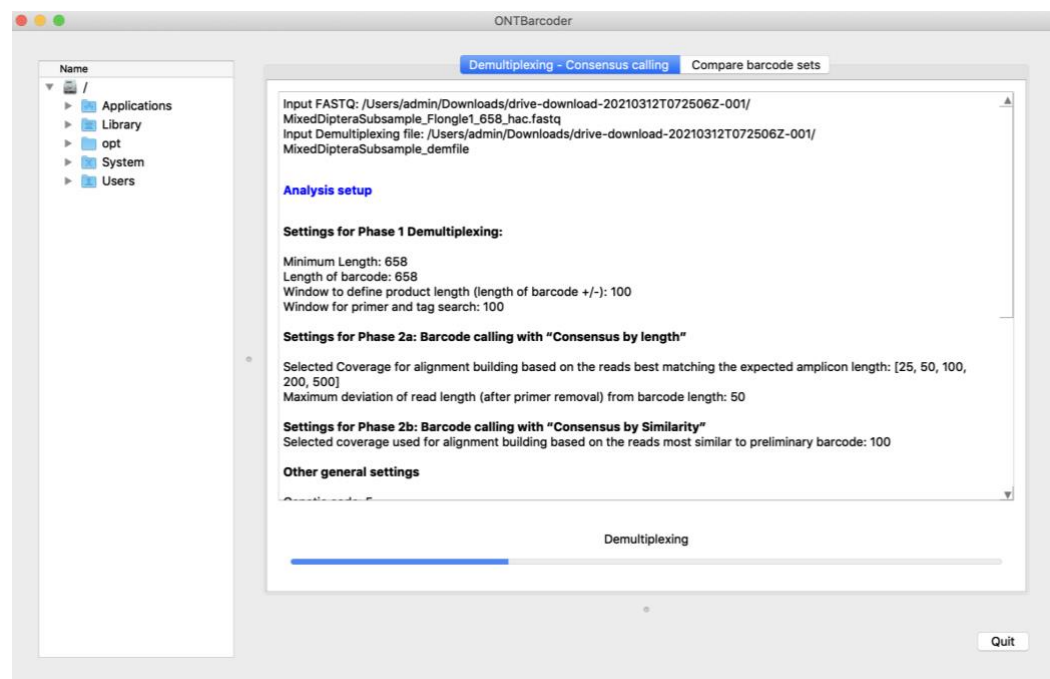
Callout boxes:

- Primers & tags are searched within the first/last specified number of bp of a read (points to 'Window for primer and tag search')
- Window size for determining product length (explen+primer/tag length)+- window size determine product length (points to 'Window to define product length')
- Main consensus calling frequency, i.e. at least 30% of the bases should be occupied by a single nucleotide for a consensus to be called (points to 'Main consensus calling frequency')
- Enables/disables **Consensus by barcode comparisons**, i.e. barcode fixing that fills/eliminates gaps in barcodes by comparison to other barcodes in the run (points to 'Consensus by barcode comparisons')
- Enables/disables **Consensus by length**: Step that uses reads closest to the known length of the barcode are used for consensus calling. (points to 'Consensus by length')
- Enables/disables **Consensus by similarity**: Step that improves consensus in previous step by taking reads most similar to it. (points to 'Consensus by similarity')
- Other consensus calling frequencies to test (points to 'Range of frequencies to assess')
- Specifies the steps for the range. In current settings 0.2,0.25,0.3,0.35,0.4,0.45,0.5 are tested (points to 'to be examined at step size of')

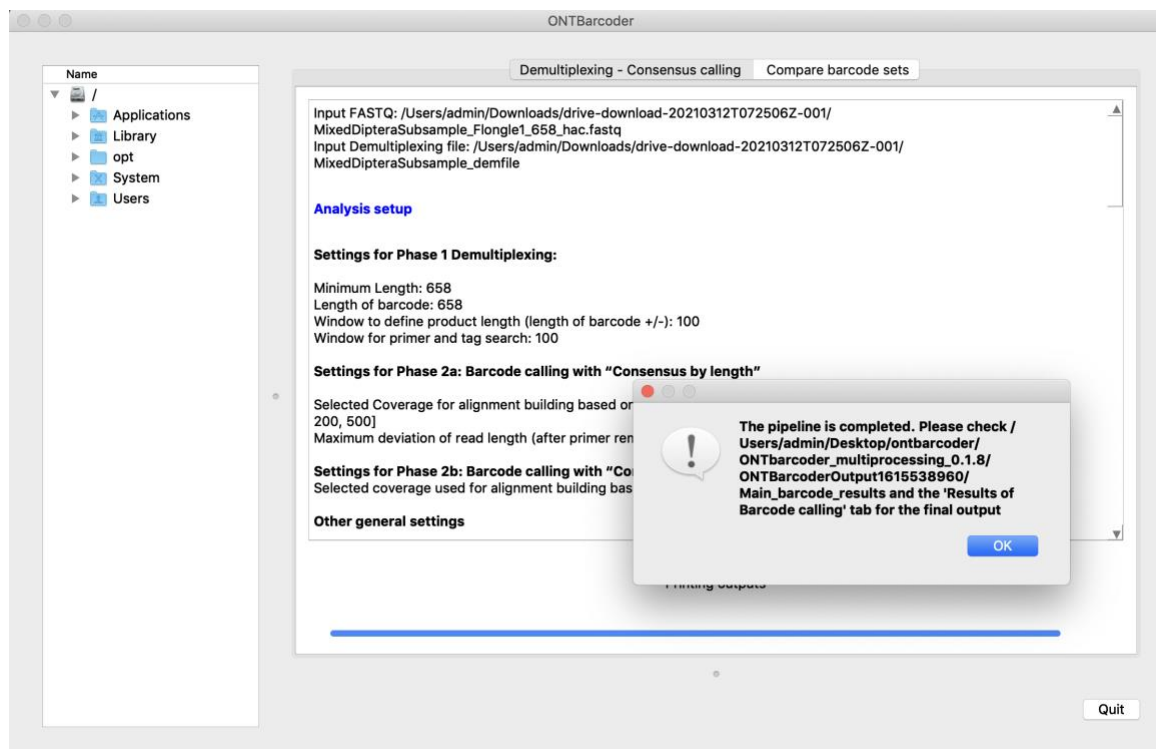
5. Select an output folder. If you select “default”, it will create a folder within the directory containing the software with a unique name starting with ONTbarcoderOutput. If you choose “Select Directory” ensure that the directory name has no empty space and is empty. Generally avoid complex characters in your names (spaces, brackets, slashes, colons)



6. The analysis starts, a progress bar appears, and results are generated in respective folders as the analysis advances.



7. The analysis is complete and all output files have been created.



8. Results Table: The table reports that 242 "Filtered" barcodes were obtained of which 192 are QC-Compliant. This means they have no ambiguities, are translatable, have expected barcode length and are not having any gaps in an internal Multiple Sequence Alignment (MSA) check. The filtered sets include addition 50 barcodes that have <1% ambiguous bases, are translatable, have expected length, and may have been corrected for up to 5 indels. The table can be copied but it also stored in the "[Outputtable.csv](#)" file in the output folder.

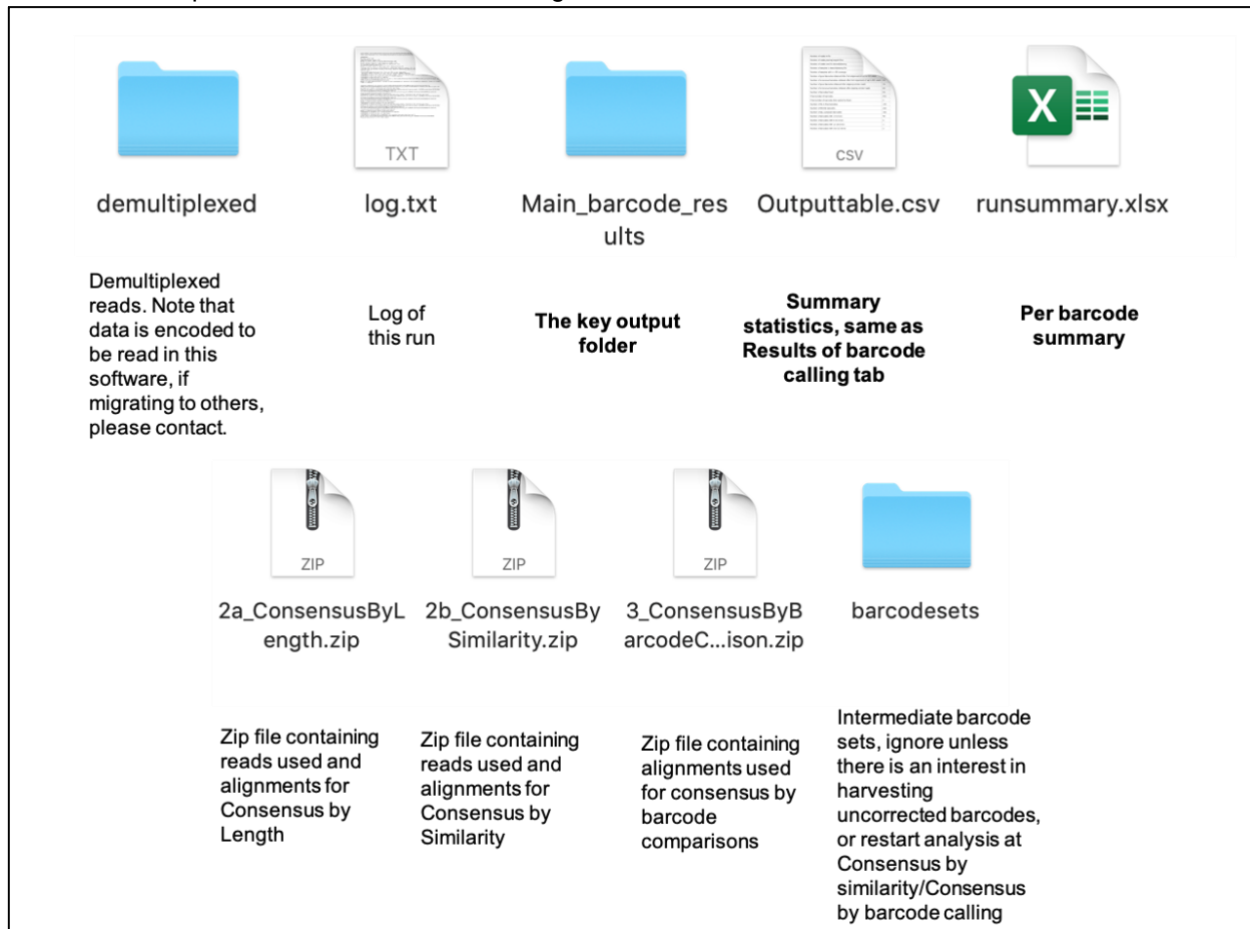
ONTBarcoder

Demultiplexing - Consensus calling Compare barcode sets Results of Barcode calling

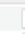


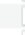








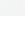


	1	2
1	Number of reads in file	294896
2	Number of reads passing length filter	222189
3	Number of reads used for demultiplexing	190952
4	Number of samples in demultiplexing file	264
5	Number of samples with >=5X coverage	250
6	Number of good barcodes obtained after first alignment of up to 200 reads	167
7	Number of erroneous barcodes obtained after first alignment of up to 200 r...	69
8	Number of good barcodes obtained after aligning similar reads	15
9	Number of erroneous barcodes obtained after aligning similar reads	68
10	Number of barcodes fixed	69
11	Final number of barcodes of expected length	251
12	Final number of barcodes that cannot be fixed	0
13	Number of Ns in final barcodes	176
14	Number of filtered barcodes	242
15	Number of QC_Compliant barcodes	192
16	Number of barcodes with 1-5 errors	48
17	Number of barcodes with 6-10 errors	5
18	Number of barcodes with 11-15 errors	1

Quit

9. The output folder contains the following



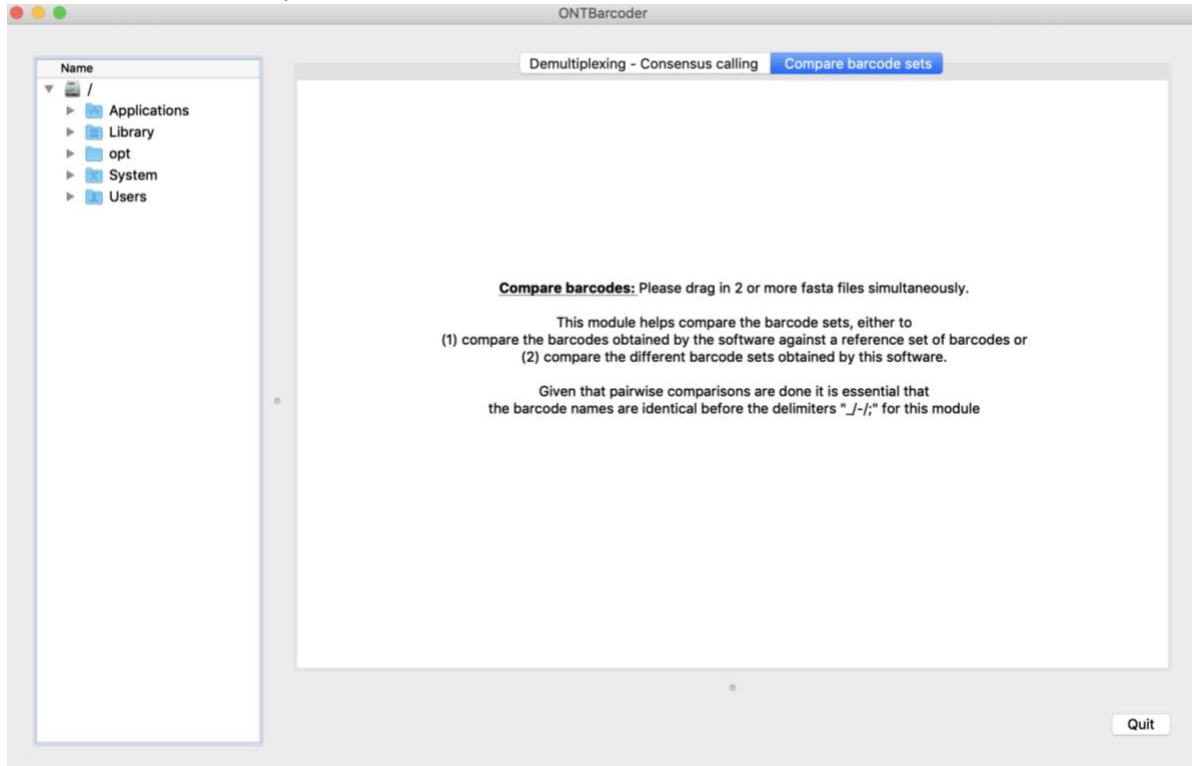
10. The Main Barcode Results folder contains the main barcode output file "**Filtered_barcodes**" fasta file. The barcodes without ambiguous bases are in the QC_Compliant_barcodes file. The zip folders contain demultiplexed reads per dataset divided into each of the categories. All barcodes should be checked via BLAST for contamination and/or verification that the barcode obtained belongs to the expected taxon.

	Allbarcodes.fasta
	Filtered_barcodes_1percamb_upto5err.fasta
	Fixed_barcodes_1to5err.fasta
	Fixed_barcodes_6to10err.fasta
	Fixed_barcodes_11to15err.fasta
	Fixed_barcodes_16to20err.fasta
	QC_Compliant_barcodes_noamb_noerr.fasta
	Remaining.fasta
	1to5errors.zip
	6to10errors.zip
	11to15errors.zip
	Filtered.zip
	Over16errors.zip
	QC_Compliant.zip
	Remaining.zip

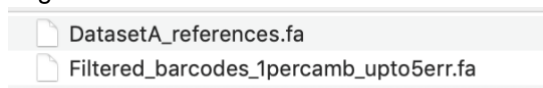
2. Compare barcodes in the “Compare barcodes sets tab”

This can be used in different modes (4 modes: single barcode file against single reference, multiple barcode files against single reference, pairwise comparisons of two barcode files, and all-vs-all comparisons of multiple barcode files), which depend on how many input files are dragged in. Here we show how to compare the output for Dataset A with reference sequences generated with Sanger sequencing.

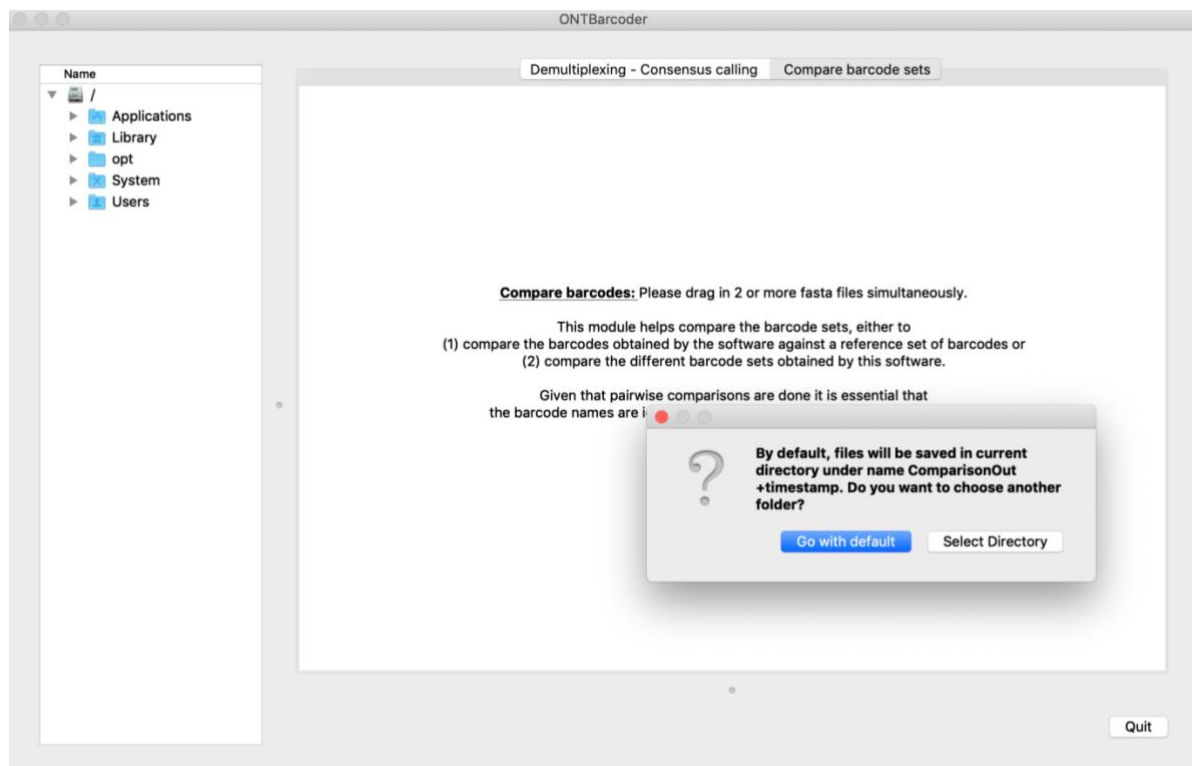
1. Switch to “Compare barcode sets” tab



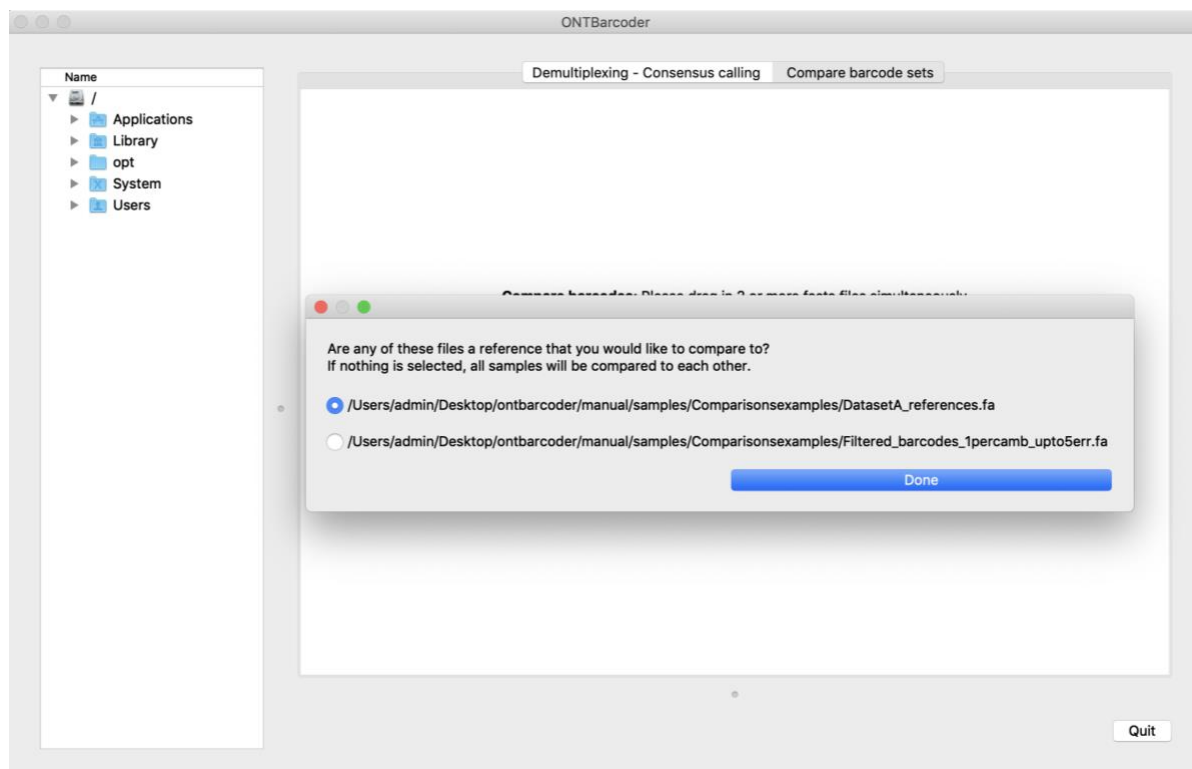
2. Drag at least two fasta files into the window. For the example, go to the [“Main barcode results”](#) folder and use the following two files:



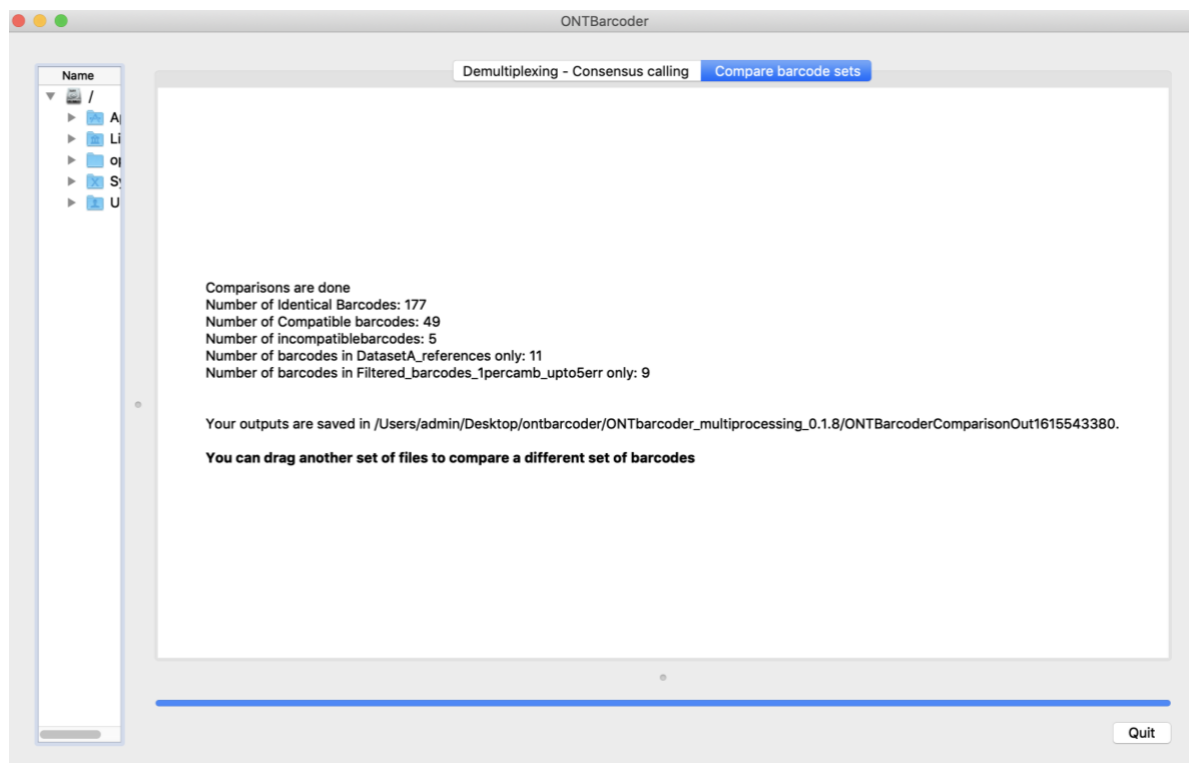
3. Select an output mode, if you select default, it will create a folder within the directory containing the software starting with ONTBarcodeRComparisonOut. If you choose “Select Directory” ensure that the directory name has no empty space and is empty. Generally avoid complex characters in your names (space, brackets, slashes, colons)



4. You can now select a set of reference barcodes if you want all other files to be compared to the barcodes in the reference set. Alternatively, you can just press “Done” to do all pairwise comparisons between the fasta files. Here, since we want the reference mode, we select the reference Sanger file and click “Done”

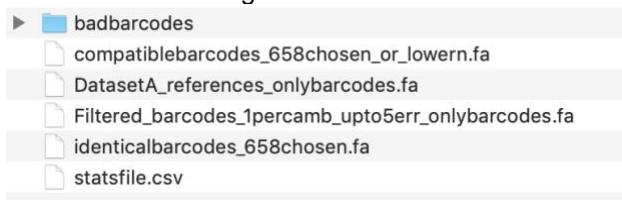


5. The output will be as shown below and the text can be copied. It is also in the output folder

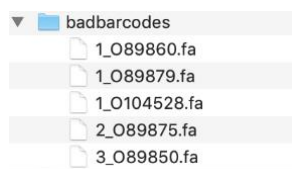


The barcodes in the sets are compared and classified into three categories: “identical” where sequences are a 100% match and lack ambiguities, “compatible” where the sequences only differ by ambiguities, and “incompatible” where the sequences differ by at least one base pair.

6. The output folder contains the following:



- “[statsfile.csv](#)” is the key summary folder having summary of the comparisons and the various incompatible barcodes with the edit distances
- Identicalbarcodes* and compatible* barcode files gives one sequence per dataset, in case the user wants to dereplicate barcodes
- The *onlybarcodes.fa files contain barcodes that cant be compared
- “badbarcodes” folder contains incompatible barcodes and the value before “_” represents edit distance.



Note that if more than 2 input files were provided, each pairwise comparison is given in a separate folder and a summary of all comparisons is then given the outermost folder

Summary

Expected input and output

1. Demultiplexing module	Input required	Results
Demultiplexing	<ul style="list-style-type: none"> FASTQ file obtained after base-calling Demultiplexing file 	<ul style="list-style-type: none"> Demultiplexed reads in “demultiplexed” folder Overall summary in “Outputtable.csv” Per barcode summary in “runsummary.xlsx”
2. Barcode calling module	Input required	Results
Barcodes are derived from reads in specimen-specific bins	<ul style="list-style-type: none"> Folder containing Demultiplexed data 	<ul style="list-style-type: none"> Barcode sets in “Main barcode results” folder Overall summary in “Outputtable.csv” Per barcode summary in “runsummary.xlsx”
Improving barcodes using “Consensus by similarity” and “Consensus by barcode comparisons”	<ul style="list-style-type: none"> Sequences in “barcodesets” folder 	<ul style="list-style-type: none"> Barcode sets available from “Main barcode results” folder Overall summary in “Outputtable.csv” Per barcode summary in “runsummary.xlsx”
3. Barcode comparison module	Input required	Results
Comparison of barcodes to references	<ul style="list-style-type: none"> One or more barcode fasta file(s) Reference fasta file <i>See format specifications</i> 	<ul style="list-style-type: none"> “statsfile.csv” that describes the overall summary and erroneous barcodes
Comparison of barcodes to each other	<ul style="list-style-type: none"> Two or more barcode fasta file(s) <i>See format specifications</i> 	<ul style="list-style-type: none"> “statsfile.csv” that describes the overall summary and erroneous barcodes

Format Specifications

File	
FASTQ file	Standard fastq, generated after basecalling with ONT software
Demultiplexing File	<p>A 5-column csv file with the following headers: SpecimenID, TagFsequence, TagRsequence, PrimerF, PrimerR</p> <p>You can only demultiplex one primer pair at a time. FASTQ files with data for multiple pairs, have to be processed sequentially.</p> <p>Please avoid unusual characters in Specimen ID (e.g. characters like “(){}[]V.,;\$” will lead to crashes)</p>
Files for comparison module	<ol style="list-style-type: none"> Specimen barcodes in different files should have identical name before the following delimiter characters (“_/-/;”) The barcodes for at least some of the specimens should be present in multiple input files
Input file for improvement (Step 4)	<p>ONTbarcoder’s pipeline can be started at different points. Demultiplexing files can be used to only carry out barcode calling. Consensus by barcode can be started with files that have the format specified in the hyperlink</p>