

# CIDR Metagenomics Hub

---

*CIDR*

*None*

## Table of contents

---

1. Introduction	3
1.1 Metagenomics network hub	3
2. Lab Resources	5
2.1 Panmetagenomics protocol	5
3. Bioinformatics	6
3.1 Clinical metagenomics bioinformatics	6
3.2 Setting up CIDR Metagenomics bioinformatics workflow	10
3.3 Starting a sequencing experiment in MinNOW	12
3.4 Running the metagenomics workflow	20
3.5 Bioinformatics - Organism query	24
3.6 mSCAPE upload tool	31
3.7 Summary report generator	34
3.8 Setting up CIDR Metagenomics bioinformatics workflow - alternative deployment	39
4. Analysis	41
4.1 Service evaluation report SOP	41
4.2 Network validation outline	42
5. FAQ	43
5.1 FAQ	43

# 1. Introduction

---

## 1.1 Metagenomics network hub

---

**The Network Hub is a resource for users of the CIDR clinical metagenomics workflow. Here, you can find SOPs, technical and FAQ/troubleshooting information regarding the implementation of metagenomics in a clinical evaluation/research setting.**

### 1.1.1 Network Sites

---

### 1.1.2 Lab protocols

---

The lab protocol is a same-day DNA/RNA extraction, host-depletion and ONT library preparation workflow for delivery of preliminary sequencing results in < 6 hours.

### 1.1.3 Informatics workflow

---

The workflow covers the end-to-end processing of respiratory samples sequencing data, delivering a metageconomic report describing the microbial communities within them. The workflow leverages ONT Nanopore sequencing at its core to produce real-time sequencing data on the GridION platform. The informatics workflow runs locally alongside the sequencing experiment, producing reports as early as 30 minutes after commencing sequencing.

### 1.1.4 Reporting framework

---

This SOP is followed to parse results from the informatics workflow for application in a clinical evaluation service setting.

## 2. Lab Resources

---

### 2.1 Panmetagenomics protocol

---

#### Note

Please reference use of this method in any presentation or publication as:

Alcolea-Medina, A., Alder, C., Snell, L.B. et al. Unified metagenomic method for rapid detection of microorganisms in clinical samples. Commun Med 4, 135 (2024). <https://doi.org/10.1038/s43856-024-00554-3>

#### Important

**Risk assessment for handling respiratory samples needs to be performed by each laboratory.**

**The tools and documentation described here and on the CIDR GitHub are not validated for diagnostic use and are for research and evaluation purposes only.**

#### 2.1.1 Version 1.2

---

#### 2.1.2 Version 1

---

## 3. Bioinformatics

---

### 3.1 Clinical metagenomics bioinformatics

---

#### Recent updates

**02/07/24** - Use of 'hospital number' have been changed to 'anonymised identifier'. Matagenomics and mSCAPE GUIs have been updated to reflect this. The mSCAPE Launcher's hospital number encryption function has been removed.

#### 3.1.1 Introduction

The principal output of the CIDR Metagenomics workflow is a PDF report listing organisms with detectable nucleic acids (RNA/DNA) and some additional information on AMR associated sequence data. The solution packages two applications - CIDR Metagenomics Workflow and [Organism Query](#) alongside a few scripts to help manage and analyse outputs. The Metagenomics Workflow runs ontop of MinKNOW, analysing sequencing data in real time producing easily digested report. [Organism Query](#) can be used to scrutinise classifications contained within a report. It leverages the full NCBI nt and RefSeq databases producing a report similar to NCBI BLAST in ~15 minutes. The Organism Query report is designed to provide the user with appropriate information to scrutinise a significant taxonomic classification.

#### The key stages of performing the bioinformatics workflow are as follow:

1. [Install workflow](#) (only on first use).
2. [Start MinKNOW sequencing experiment](#).
3. [Start the Metagenomics workflow](#).

#### Optional:

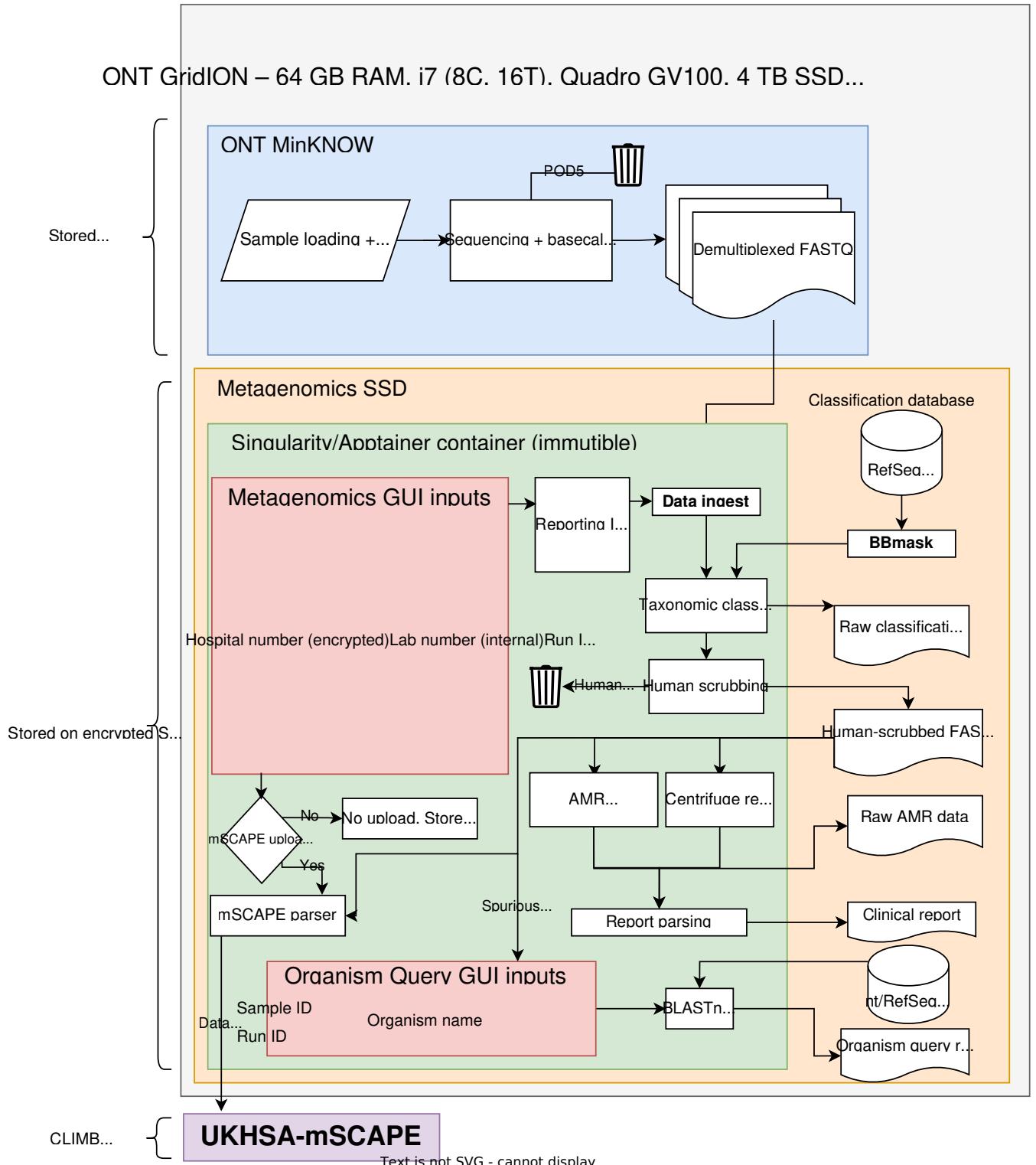
1. [Query a classification](#).
2. [Upload data to mSCAPE](#).
3. Generate a summary spreadsheet.

#### mSCAPE

Users can opt in to [mSCAPE](#) on an per-experiment basis for an automatic upload of sequencing data to UKHSA mSCAPE.

## Technical facets

After loading a metagenomic library on to an ONT sequencing device and launching the sequencing experiment in ONT MinKNOW the pipeline is initialised by the user through the Metagenomics Launcher graphical user interface (GUI). The software periodically ingests base called FASTQ data from the GridION `/data/` directory at set intervals - 0.5, 1, 2, 16 and 24 hours. At each interval, the pipeline performs human scrubbing, taxonomic classification, AMR detection and MLST which is then consolidated in to a PDF reports which are saved in the `/media/grid/metagenomics/results/` directory. The diagram below illustrates further details of how the pipeline works:



### taxonomic classification

At its core, the pipeline leverages a taxonomic classification tool called [Centrifuge](#). There are numerous alternatives, the most common being Kraken. We chose centrifuge namely because of its smaller memory footprint and existing deployment in ONT's WIMP. Each read in the raw data is aligned against an index, and is assigned a

confidence score and a taxonomy. Our index is an optimised database of curated eukaryotic, prokaryotic, and viral reference sequences formed primarily from [NCBI RefSeq](#) and [FDA-ARGOS](#) databases. This is provided on the SSD sent to each site.

The index was built from the following [sequences](#) - this map file contains the accession number for each sequence in the database and the accompanying taxa ID. This file is required to assemble the Centrifuge index alongside the sequences in FASTA format. Before building the index, the sequences were masked using [BBmask](#). This tool is applied primarily to prevent false-positive matches in highly-conserved or low-complexity regions of genomes.

### 3.1.2 Related code snippets

---

#### Masking a FASTA database using BBmask

```
bbmask.sh in=unmasked.fasta out=masked.fasta entropy=0.7 -Xmx80g maskrepeats=t
```

#### Building a centrifuge index -

```
# --bmax needs tuning based on available memory. centrifuge-build -p 10 --conversion-table accession2taxid.map --taxonomy-tree ./taxdump/nodes.dmp --name-
```

## 3.2 Setting up CIDR Metagenomics bioinformatics workflow

---

### 3.2.1 Overview

Each Network site will receive an ONT GridION sequencing platform and an external SSD containing the software and databases required for analysing metagenomic datasets. The software has been designed such that it will be easy for anybody to set up and use. Follow the instructions below to install the bioinformatics workflow.

### 3.2.2 Install instructions

1. Insert the USB SSD in to one of the blue USB ports at the rear of the GridION. Try to place the disk away from the warm exhaust as this may lead to overheating.
2. After logging in to the GridION Ubuntu operating system, modify the file browser setting by following the video below. This is to enable the running of scripts without using the terminal.
1. Using the file browser, on the taskbar on the left side of the screen, navigate to the **metagenomics** disk, which can be found in the navigation pane inside the file browser.

#### Info

As a security feature, the removable SSD has been encrypted. Enter the encryption key provided to you and confirm that you'd like the key remembered.

1. Navigate to the `metagenomics` disk in the file explorer and double-click `launch_installer.sh`, selecting to 'Run in terminal'. When prompted to do so, type the password for the GridION device (not the encryption key). See below for a video guide.

#### Info

As you type, no lettering or symbols will appear. This is normal. If you mistype, press enter and try again.

There may be additional outputs in your terminal window compared to the video.

1. Some icons should appear on the desktop linking to each app. You will need right click on the icons and select `Allow launching` before continuing with the next stage.

Success!

We have now installed the CIDR metagenomics workflow. The next step will be to run through a control dataset to test the workflow has run sucessfully.

### 3.2.3 Install validation

Included with the software is a small dataset based on the Zymo community standard. In this step we will validate the function of the workflow with this dataset and generate a report.

1. Double click the Metagenomics Launcher icon on the desktop.
2. Fill out the fields, as indicated in the video below. More information on how to fill the fields and run the launcher can be found in the [Starting the metagenomics workflow](#) section.
1. Wait ~10 minutes for the workflow to complete. Open up the PDF report which can be found in the `/reports` folder on the metagenomics disk in a folder corresponding to the name of the sample provided in the launcher eg. `gstt_control_1`. See video below for further information.
1. Inspect the `/metagenomics/reports/validation_sample/` PDF report at the **0.5 hr timepoint**, it should match the CIDR validation report provided [here](#).

Success!

We have now tested the CIDR metagenomics workflow. The next step will be to run a sequencing experiment, running the workflow in real-time.

## 3.3 Starting a sequencing experiment in MinNOW

---

### 3.3.1 Introduction

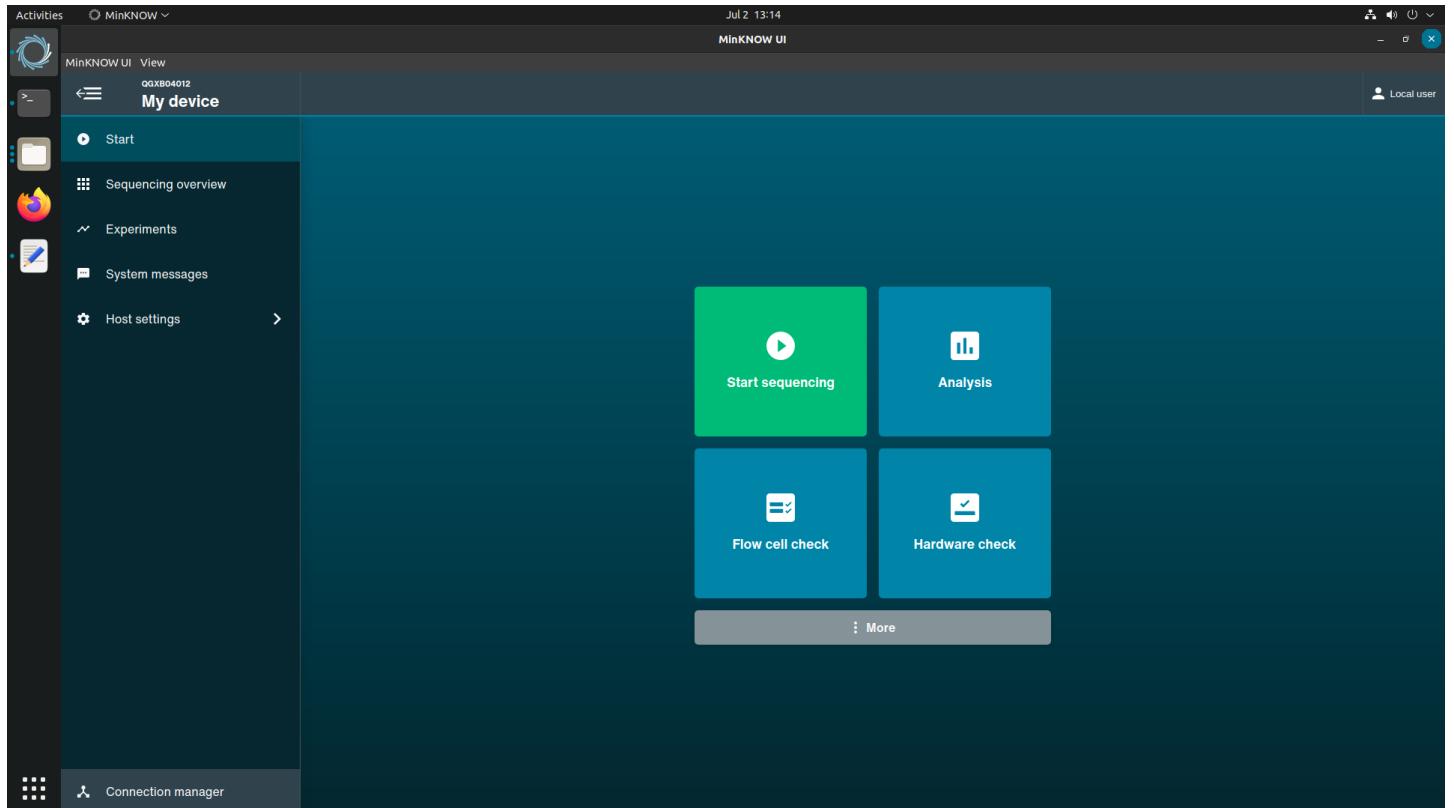
This section documents user interactions with the ONT MinNOW interface. The MinNOW software is used to initialise and control sequencing experiments on the GridION device. The CIDR Metagenomics Workflow is launched after having completed this section.

Before completing any protocols, users should check flow cells sent to them by ONT are above the warranty pore count. The pore check should be run:

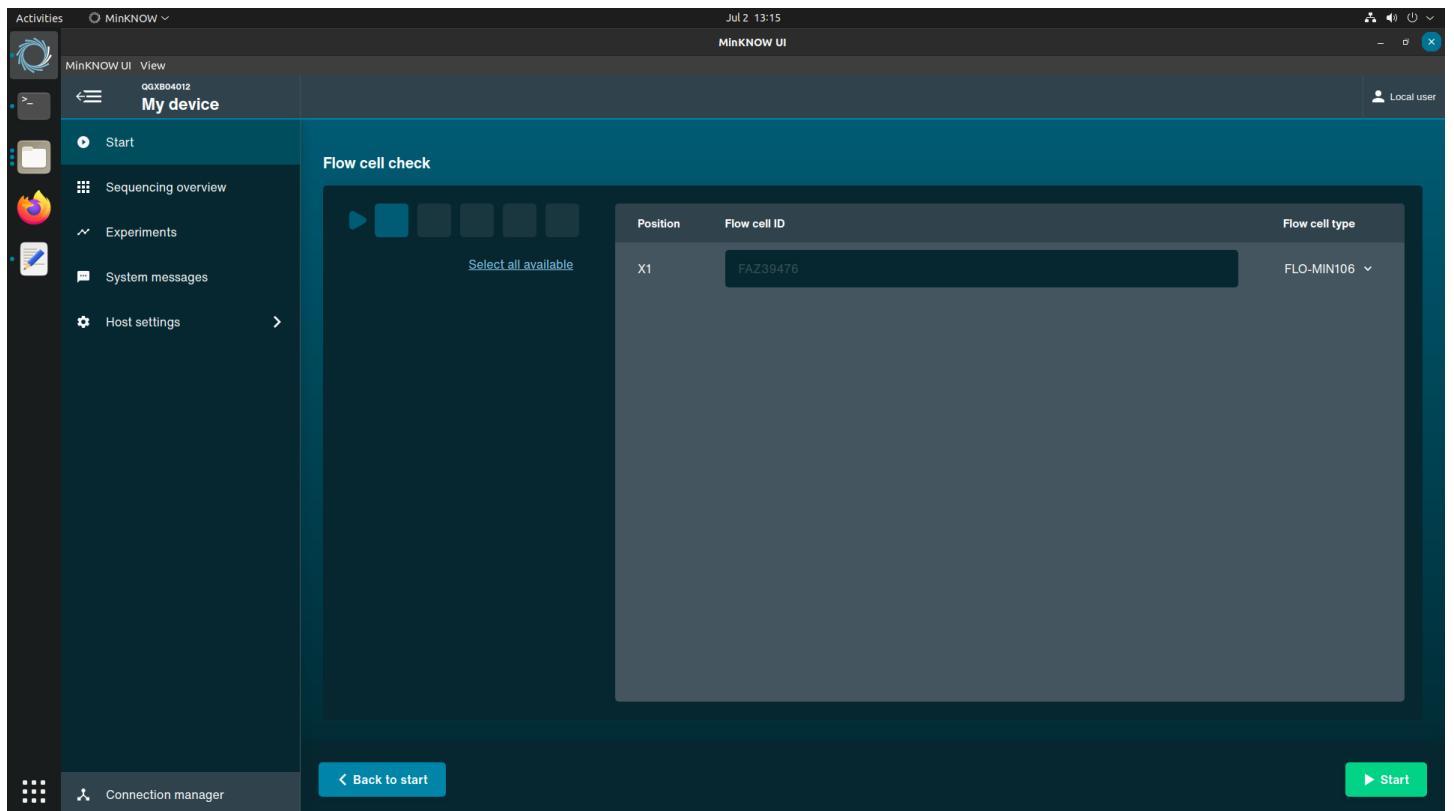
1. As soon as you have received a new batch from ONT.
2. Immediately before starting a sequencing experiment.

### 3.3.2 Running a flow cell check

Oxford Nanopore Technologies will replace any flow cell that falls below the warranty number of active pores within three months of purchase, provided that you report the results within two days of performing the flow cell check and you have followed the storage recommendations. A MinION flow cell (used also in the GridION) should have 800 pores.



1. Select **Flow Cell Check** from the MinKNOW Start screen.
2. Indicate the corresponding sequencing positions you'd like to check by selecting the square icons below the 'Flow Cell Check' title. (See image below)
3. Click on the green Start button and wait for the flow cell check to complete.
4. If the pore count is < 800 and the flow cell is still in warrant, contact ONT for a replacement within two days of completing the check.



### 3.3.3 Starting a sequencing experiment on MinNOW

1. From the Start screen in MinNOW select 'Start Sequencing'.
2. Select the position occupied by the flow cell loaded for the sequencing experiment, enter the Experiment and Sample IDs and select 'Continue to kit selection' at the bottom of the screen.

#### Note

You should enter a new and unique Experiment ID and Sample ID for each new sequencing library.

A summary of the configuration parameters is shown at the bottom of this section.

1. Select the RPB-004 library preparation kit from the Kit selection screen. Click 'Continue to run options' at the bottom of the window.

Kit selection

Sample type	PCR-free	Multiplexing	Control	
DNA	PCR	Yes	<input type="checkbox"/>	Reset filters
Ligation Sequencing Kit SQK-LSK109	PCR Barcoding Kit SQK-PBK004	16S Barcoding Kit (BC1-24) SQK-16S024	Ligation Sequencing Kit (48 reactions) SQK-LSK109-XL	
Ligation Sequencing Kit SQK-LSK110	Ligation Sequencing Kit XL SQK-LSK110-XL	Rapid PCR Barcoding Kit SQK-RPB004	VoITRAX PCR Tiling 1-12 COVID-19 VSK-PTC001	<input checked="" type="checkbox"/>
VoITRAX RT-PCR Sequencing Kit 1-12 VSK-VPS001	16S Barcoding Kit SQK-RAB204			

< Back to position selection      Continue to run options >      Skip to final review >

1. In the run options screen, set the sequencing experiment to last for 24 hours. Leave the read length at 200 bp and the other settings as default and select 'Continue to analysis' at the bottom of the window.

The screenshot shows the MinKNOW UI interface. The top navigation bar includes 'Activities' and 'MinKNOW' dropdowns, a date and time indicator ('Jul 2 13:19'), and a user icon ('Local user'). The main header 'MinKNOW UI' and 'My device' are displayed. On the left, a sidebar lists 'Start', 'Sequencing overview', 'Experiments', 'System messages', and 'Host settings'. The main content area is titled 'Run options' and contains the following sections:

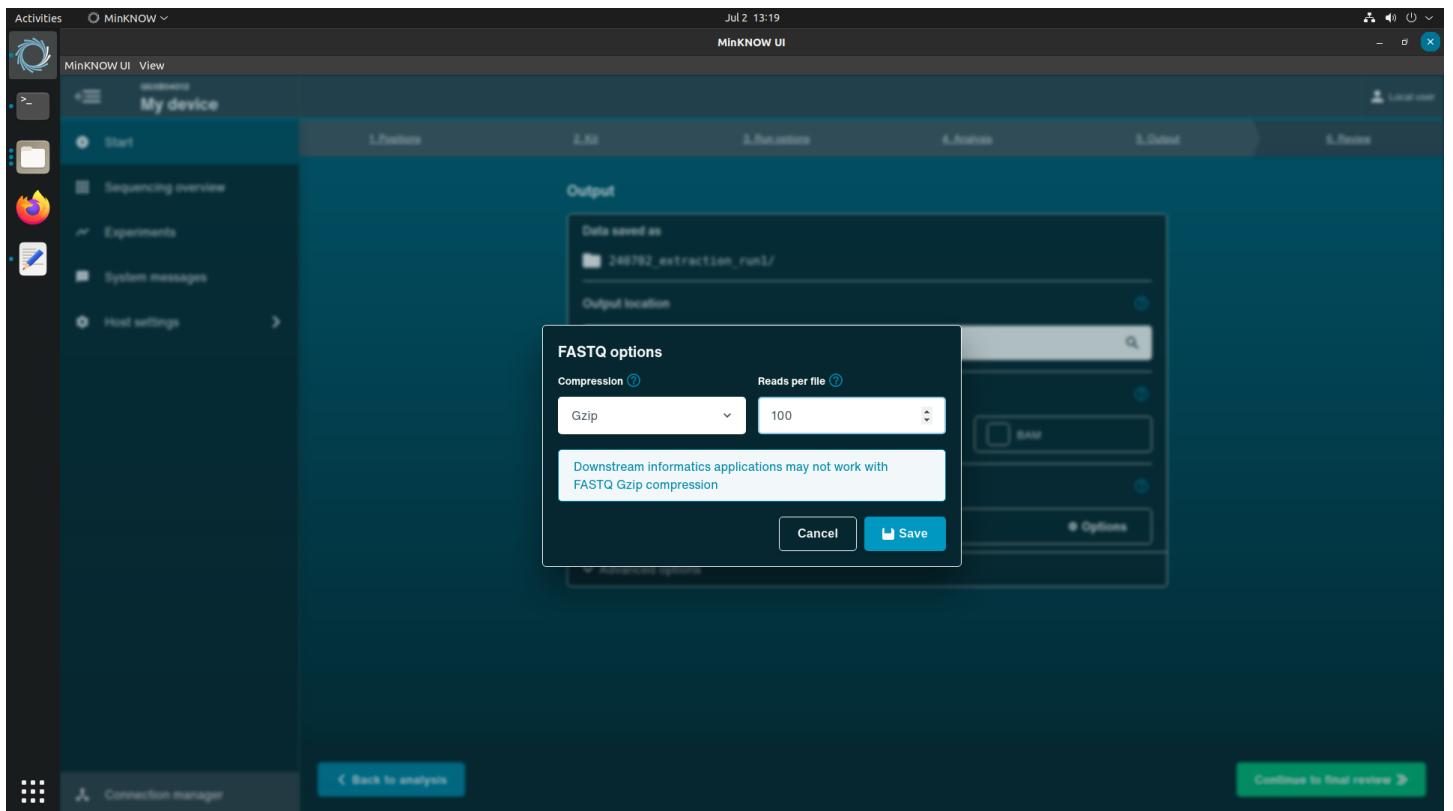
- Run duration:** A slider set to '24 hours'.
- Minimum read length:** A slider set to '200 bp'.
- Adaptive sampling:** Options include 'Enrich or deplete sequences' and 'Barcode balancing' (marked as 'Beta').
- Advanced options:** A collapsed section.

At the bottom, there are buttons for 'Back to kit selection', 'Continue to analysis >', and 'Skip to final review >'. The overall theme is dark blue.

1. On the Analysis window, under barcoding, select 'Edit options'. In the popup window, select 'Barcode at both ends' and 'Mid-read barcode filtering'. Select the 'Continue to output' button at the bottom of the window.



1. On the Output window, deselect the FAST5 option. Where the FASTQ checkbox is selected, click the gear icon and set 'Reads per file' to 100. Continue to final review.



1. Check the parameters below match what is indicated on-screen.

Parameter	Value
<b>Selected Kit</b>	SQK-RPB004
<b>Run length</b>	24 hours
<b>Minimum read length</b>	200 bp
<b>Adaptive sampling</b>	Off
<b>Basecalling</b>	On (High accuracy basecalling)
<b>Barcode</b>	On
<b>Require both ends</b>	On
<b>Detect mid-read barcodes</b>	On
<b>Alignment</b>	Off
<b>Location</b>	/data
<b>FAST5</b>	Off
<b>FASTQ</b>	On (Gzip, 100 reads per file)
<b>Read filtering</b>	Qscore:9 Readlength: unfiltered, Read splitting: Disabled

1. Start the sequencing experiment. After the flow cell reaches temperature, navigate to the barcodes screen to verify data has been output.

Success!

After reads start to appear on the barcoding screen you can advance to [Starting the metagenomics workflow](#).

## 3.4 Running the metagenomics workflow

### Before starting

1. The CIDR metagenomics workflow must be started during a sequencing experiment or after a sequencing experiment has completed. The pipeline must not be activated before a sequencing experiment has started in MinNOW and has **started producing reads** (See MinNOW setup - Lab Protocol).
2. Ensure the SSD is inserted in to one of the rear USB 3.1 ports, has been mounted and the encryption key has been entered successfully. Test the disk has been mounted by navigating to it in the Ubuntu file explorer.

### Starting a run

1. Double click the **Metagenomics Launcher** icon on the GridION desktop, the CIDR Metagenomics Launcher should appear alongside a terminal window.



## Known issues

The `'geocryptfs error not found...'` error can be ignored as it is not essential to the workflow.

If a sample is repeated, append the Lab ID accordingly (\_2) - eg. 123mre123456\_2

1. Select the number of samples to be analysed from the dropdown.
2. You can choose to initiate the launcher using **one** of the below methods:

- Fill out the fields on the form for each sample to be analysed.
- Loading a pre-existing TSV - [see example](#).

## Field descriptions:

Field	Description
<b>MinKNOW experiment ID</b>	The exact name matching the experiment name on MinKNOW entered by the user when initiating a sequencing run. This is populated automatically from the <code>/data</code> directory.
<b>KinKNOW sample ID</b>	The exact name matching the Sample name on MinKNOW entered by the user when initiating a sequencing run. This is populated automatically from the <code>/data/{experiment_id}/</code> directory.
<b>ONT barcode</b>	The ONT library index/barcode used. Green colour indicates the barcode directory has been validated.
<b>Lab/Sample ID</b>	The unique lab accession number for the sample. This data is encrypted before transmission. <b>If repeating a sample, append with _n</b>
<b>Anonymised identifier</b>	An anonymised identifier linked to the sample hospital number.
<b>Collection date</b>	Collection data of the sample.
<b>Sample Class</b>	The class of sample loaded.
<b>Sample type</b>	The methodology used to collect the sample.
<b>Operator</b>	Arbitrary identifier of the user operating the sequencer.

## Note

- Option 1 will generate a sample sheet stored in the `metagenomics/sample_sheets` directory. This can be reused if a repeat run is required - or quick edits need to be made to a set of samples without having to fill out the fields again.
- Filling the 'filename suffix' field will save the sample sheet with an appended string of your choosing to help identify your run's metadata in the 'sample\_sheets' folder.

1. With the metadata form filled, select the run parameter check boxes.

Parameter	Description
<b>Force overwrite</b>	The exact name matching the experiment name on MinKNOW entered by the user when initiating a sequencing run. This is populated automatically from the <code>/data</code> directory.
<b>mSCAPE prompt</b>	After the sequencing and analysis run has completed, open the mSCAPE uploader for user input. No data is uploaded without par-sample expressed authorisation.
<b>Disable data ingest sleep</b>	<b>Not for real-time analysis!</b> Analyse all data immediately - do not wait for it to be generated by the sequencer.

## Known issues

You should wait to launch the pipeline after the sequencer has reported producing reads in MinKNOW, the workflow will display errors in red if no reads have been found.

The NTC will exhibit the same 'error' behavior as no reads are present in the corresponding barcode folder. We are working on functionality to circumvent this.

You can stop the analysis or close the Launcher window at any point by closing the terminal window. The terminal window can be closed using the  in the top right corner.

1. Click on `Launch pipeline` and click `OK` to start analysis.
2. After a minute, the terminal window accompanying the workflow launcher should start displaying log outputs from the workflow. See below for an example.
1. ~40 minutes after launching the sequencing experiment alongside the metagenomics workflow, the first reports will be available in `/media/grid/metagenomics/reports/{sample_name}/{timepoint}`. See below for a guide on how to access this.

Success!

We have now run the CIDR metagenomics workflow. The workflow will run for ~24 hours generating PDF reports for 0.5, 1, 2, 16, 24 hour time-points.

## 3.5 Bioinformatics - Organism query

---

The Organism Query tool is designed to help scrutinise taxonomic classification outputs from the CIDR metagenomics workflow. It uses a local (offline) version of NCBI BLASTn, with the full NCBI nt database to produce a report, similar to that found on the NCBI BLAST website, providing the user with a second opinion on classifications.

### 3.5.1 Technical information

- The tool searches the classified reads for an organism indicated by the user. Selecting a subset of 10 (max) reads assigned to that taxa.
- The reads are extracted from the FASTQ file stored in the workflow `results` folder and BLASTs them against the prescribed database.
- The results are parsed in to a HTML report with an interactive plot designed to help the user explore the different metrics of alignment.
- A full BLAST alignment report is available at the bottom of the HTML report.
- The report is stored in the `reports/{sample_id}/organism_query_XXX` folder, with the other PDF reports from the metagenomics run.
- In the report folder is the BLAST HTML report and the subsetted FASTQ and FASTA reads.

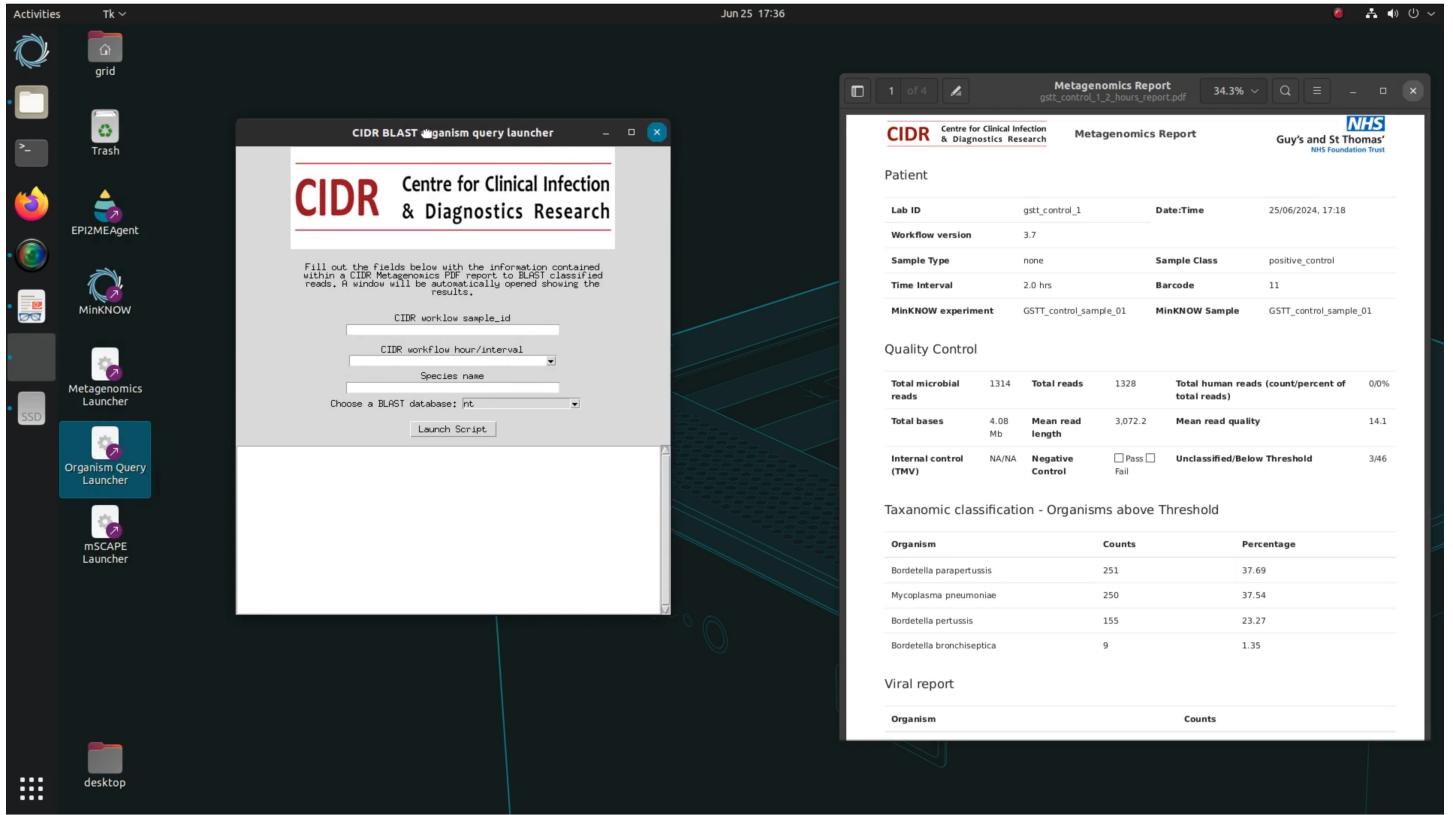
### Info

The most thorough analysis is performed using the default nt database. This usually takes ~30 minutes to generate a report.

### Known issues

The interactive plot is only available immediately after running the tool, when the original Launcher window is open. Loading the query report after closing will allow you to view the BLAST alignments but not the plot section, which will display an 'Internal Server Error'.

Running queries simultaneously is currently not supported as the interactive plot function can run a single session. If two reports are opened, the incorrect plot may display on the report. This does not effect the full BLAST analysis, only the plot displayed at that time.



### 3.5.2 Launching organism query

1. Load the PDF report from the run you'd like to query a classification from.
2. Click on the Organism Query launcher icon on the GridION desktop.
3. Fill out the fields as indicated in the video below.

Parameter	Description
<b>SampleID</b>	The Lab/Sample ID assigned to the sample when launching the metagenomics workflow.
<b>Workflow hour/interval</b>	The timepoint, corresponding to the dataset you'd like query.
<b>Species name</b>	The name of the species to be queried eg. Aspergillus fumigatus

1. Click on the **Launch script** button to start the query workflow. A Firefox browser window will appear after the workflow has finished. You can reopen the report from `reports/{sample_id}/organism_query_XXX` on the metagenomics SSD.

### 3.5.3 Interpreting results

---

Interpreting the results can be subjective. The interactive plot has been designed to help guide decision making. Using a combination of the alignment length (relative to the read length), the identity score and the e-score can be an informative approach.

- Query Coverage: Percent of the query sequence length that is included in alignments against the sequence match.
- E-value: Indicates the number of hits or alignments that are expected to be seen by random chance with the same score or better. The lower the E-value, the more significant the alignment (the closer to 0, the better). E-value is the default metric used to sort the Descriptions table. Click here for a discussion of E-value thresholds.
- Percent Identity: Percent of nucleotides or amino acids that are identical between the aligned query and database sequences. A query sequence can share low percent identity with a sequence and still be a significant hit. It is essential to take the E-value into account and look for similarity between conserved regions (this will be more evident at the amino acid level).

## Interactive plot

In a situation where there is either no confident match or no convergence across the read subset on an organism, or alignments indicating a number of similar, closely related organisms, the plot appears disordered in a single cluster.



Situations where a convergent set of alignments are indicated results in a plot with two clusters, usually with one clear homogenous taxa. In this case, the Tobacco Mosaic virus.



## BLAST alignments

Scrolling to the bottom of the report is a conventional BLAST alignment. We recommend looking at each read, the **Length** (query readlength), the **Score** and **Identities** values.

- If only a small alignment (80 bp) has been made from a long 2000 bp read, this is not likely to be a robust estimation.
- Check the alignment visualisation for long repeats, regions of low complexity etc. These regions often confound BLAST and taxonomic classification tools.

Read number

0	1	2	3	4	5	6	7	8	9
---	---	---	---	---	---	---	---	---	---

Query= d056f53c-29e1-4b5a-b19a-46478034fbcc

Length=3020

Sequences producing significant alignments:

	Score (Bits)	E Value
<a href="#">AP019841.1</a> Leptotrichia wadei JMUB3936 DNA, complete genome	4010	0.0
<a href="#">AP019827.1</a> Leptotrichia shahii JCM16776 DNA, complete genome	2719	0.0
<a href="#">AP019829.2</a> Leptotrichia wadei JCM16777 DNA, complete genome	2357	0.0
<a href="#">AP019834.1</a> Leptotrichia wadei JMUB3933 DNA, complete genome	2351	0.0
<a href="#">AP019835.1</a> Leptotrichia wadei JMUB3934 DNA, complete genome	2340	0.0

>[AP019841.1](#) Leptotrichia wadei JMUB3936 DNA, complete genome  
Length=2335974

Score = 4010 bits (2171), Expect = 0.0  
Identities = 2693/2919 (92%), Gaps = 140/2919 (5%)  
Strand=Plus/Minus

Query	96	TAGATGAAAACGGAAATGTACCAATCGCAAGGACGTGCCCTAATGGCAGATGCAATAGCA	155
Sbjct	898962	TAGATGAAAACGGAAATGTACCAATG-TAGGACGTGCCCTAATGGCAGATGCAATAGCA	898904
Query	156	ACTACGGCTGGCGCAGCAC--CCA--TTCAACGGTTACAGCTTATGTGGAAAGCTAAC-	210
Sbjct	898903	ACTACAGCTGGCGCAGCACTTGGAGTTCAACAGTTACAACTTATGTGGAAAGTTCAACA	898844
Query	211	GG-G--GTCGCGGGTGGAAAGAACTGGATGGACCTTCATCACAAACAGGAGTTTATTCTA	267
Sbjct	898843	GGAGTTATCGCAGGCGGAAGAACTGGATGGACAGCCATCACAAACAGGAGTTTATTCTA	898784
Query	268	ATATCAATGTTTCTCACACCATAATTATTCGATACCAGGATGTGCCACAGCTCCAGC	327
Sbjct	898783	ATATCAATGTTTCTCAC-CAATATTATTCAAATACCAGGATGTGCCACAGCTCCAGC	898725
Query	328	-TGGGCCA-GGTAGTTATTAATGCTAAAGTTCACTCAAAAAACACTG-GTTGCATGATG	384
Sbjct	898724	CTTAATTACGTTGGTTATTAATGCTAA-GTTCAAGTAAAAACATAGATTGATGATG	898666
Querv	385	TCTTGGAAAGGTGTTCCATCATTTATCACAACTACAATGGCTTAACCTATAGCATCG	444

## 3.6 mSCAPE upload tool

### 3.6.1 Introduction

With on-boarding to the mSCAPE programme, users are encouraged to upload samples to mSCAPE using the mSCAPE Launcher. This tool incorporates sample metadata inputted when using the Metagenomics Launcher, with some additional user inputs. The tool then packages human-scrubbed FASTQ data, producing metadata outputs in the required format for an mSCAPE submission.

### 3.6.2 Installation

The mSCAPE on-boarding team will provide all of the credentials required to make a test submission. These credentials should be stored in `~/.aws/` on the host machine. Contact the bioinformatics lead at GSTT for activation of the mSCAPE Uploader.

### 3.6.3 Uploading samples

Following the successful completion of the CIDR Metagenomics Workflow, users can double click the 'mSCAPE Launcher' on the desktop. The tool reads the sample information from the sample sheet saved when a workflow run is started and parses it in to the correct format ready for upload.

#### Note

Using the 'filename suffix' field on the Metagenomics Launcher appends the sample sheet filename with a string of your choosing, making it easier to find during audits or mSCAPE uploads.

*Check out the [video at the bottom of the page](#) for a visual guide on how to run the uploader*

1. From the 'Dropdown Options' section, select the parameters appropriate for the samples to be uploaded. Some addition fields for mSCAPE are inferred from the sample sheet, the data in the dropdowns or the sequencing reads themselves. See the table below for more details.

## Note

For more details on how to fill the metadata fields associate with the samples, for example 'Study description', [see the additional information section](#) at the bottom of the page.

Parameter	Description
<b>StudySite</b>	The RMg Network site the sequencing took place.
<b>Extraction Method</b>	the nucleic acid extraction methodology used.
<b>Spike-in</b>	The spike-in control used.
<b>ISOCountry</b>	Country/nation.
<b>Sequencing Protocol</b>	The methodology used for sample preparation.
<b>Library Protocol</b>	the sequencing kit used for library preparation.
<b>Bioinformatics protocol</b>	The version of the CIDR Metagenomics bioinformatics workflow.
<b>Clinical or research</b>	mSCAPE clinical or research.
<b>Human scrubbing</b>	mSCAPE informatic Human Scrubbing protocol used.
<b>Study description</b>	Code provided by mSCAPE team.

1. Load a sample sheet by clicking the 'Load sample sheet' button at the bottom of the interface. The file browser takes you to the 'sample sheet' directory. Find the date/time corresponding to your run (or the sheet identified by the filename suffix) and open it. With this, the Main Table section should be populated with the samples from the corresponding sequencing run.
2. From the 'FASTQ Selection', select the dataset and the timepoint to be uploaded.

## Note

Sample which are not intended for upload can be checked in the 'Main Table' and removed by clicking the 'Delete Selected rows' button at the bottom of the interface.

1. Once you have selected all of the timepoints for upload, choose the 'Update DataFrame' button at the bottom of the interface.
2. Review the data in the 'Main Table' panel ensuring it is suitable for upload.
3. Select 'Upload to S3'. You can check the terminal window for outputs confirming successful upload. A status indicator is also present in the FASTQ selection window.

## Video tutorial

*Video demonstration of using mSCAPE Uploader features*

### 3.6.4 Additional information

#### Study description

Once onboarded by the CLIMB team, the site will be provided with a predetermined [study\_centre\_id] and three [study\_id's].

A study\_centre\_id is an abbreviation of the site name e.g. name of the NHS trust. Each site will have one ID.

A study\_id will be used to identify if samples are from a particular research study, or an NHS residual sample. Every NHS trust will have one Study ID related to their clinical samples, and a separate verification Study ID. This is to clearly differentiate between their verification stage and clinical service. Additional study\_ids can be added as desired to differentiate any research studies.

A list of the predetermined site and study codes can be found here: (Link to be added soon)

**See Table 1 for example of a study\_centre\_id, study\_id and their purposes for GSTT.**

Table 1. Example of a study centre ID and study ID:

study_centre_id	study_id	Purpose
GSTT	GSTT-CLI-01	Clinical service samples
GSTT	GSTT-VER-01	Verification
GSTT	GSTT- RES-01	Research – outside of network of excellence protocol

## 3.7 Summary report generator

---

### 3.7.1 Purpose

We have found it useful to be able to generate summaries of multiple runs for downstream analysis. This takes the form of a spreadsheet or CSV, with rows corresponding to each sample and columns containing information derived from the sample sheet and the metagenomic analyses across timepoints. The tool features an end-to-end GUI to select samples and build the sheet.

#### Important note

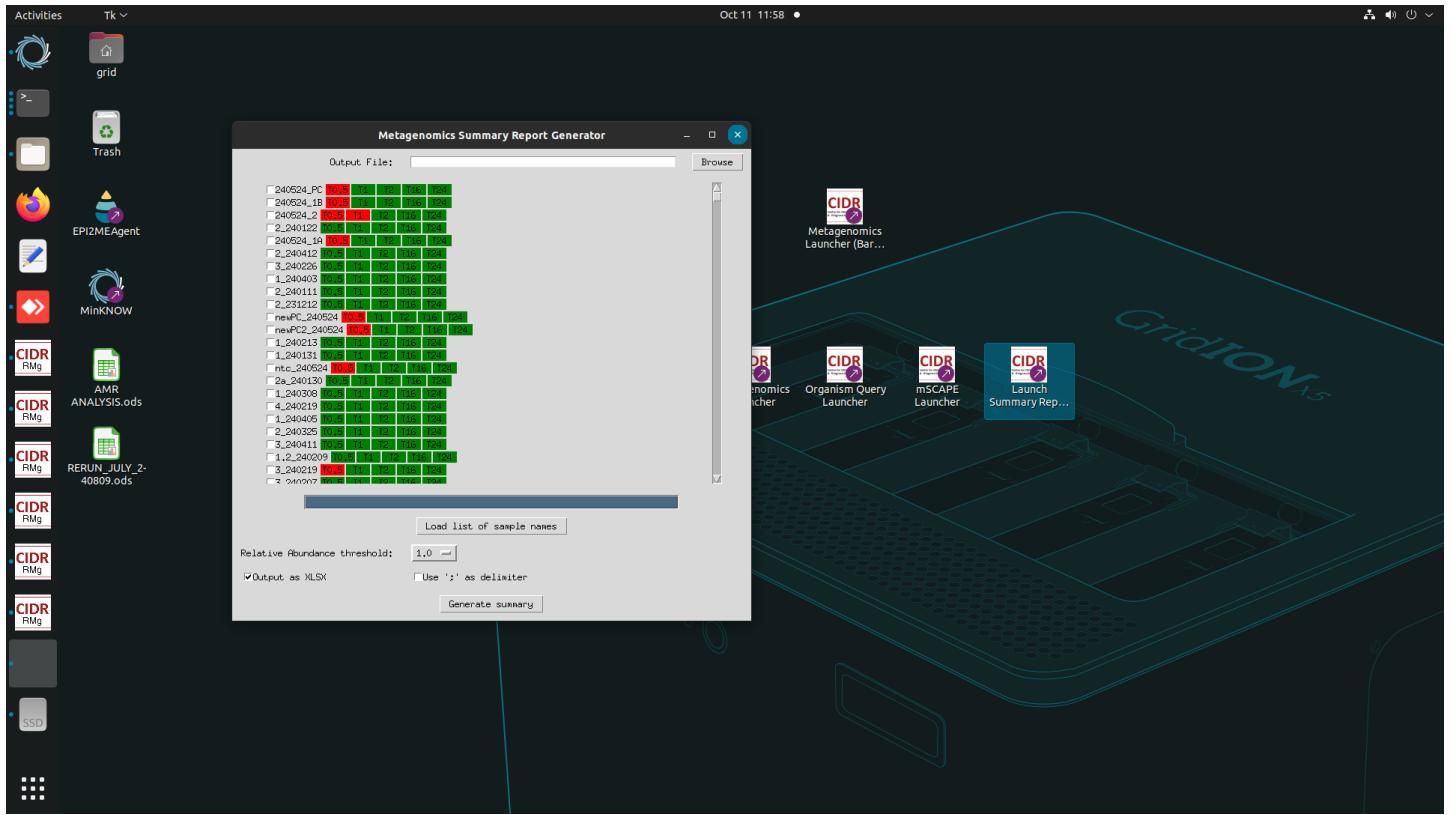
We ask that for the NHS RMg service evaluation, you set the Summary Report threshold cut-off to the lowest value (currently 0.1) so that all classifications are included.

A description of the fields featured in the spreadsheet is available [here](#).

### 3.7.2 Instructions for use

*Check out the [video at the bottom of the page](#) for an end-to-end demonstration*

1. Double click the Launch Summary Report icon on the desktop. The window should appear with a loading bar.  
Please wait until all of the samples are sourced and loaded.



## Note

In the above example, we have a list of samples, some of which have red indicators at specific time point. In both instances, it is likely that the reports were not generated because there were no reads present in the dataset. This is an especially frequent occurrence in NTC samples given that few reads should be detected in these samples.

1. The program reads all Metagenomics Workflow runs from the 'results' folder and populates a list. The list has a checkbox to include or exclude samples from the report, and a coloured box indicating the time point is present in the dataset. Select the samples for reporting using one of the two below methods.
  - a. Select samples using the check boxes on the interface.
  - b. Produce a simple list (newline delimited) of sample names, matching exactly (case sensitive) the Sample ID/Accession number used in the Metagenomics run. Save the list in the 'Sample Sheet' directory on the Metagenomics SSD. Select the 'Load list of sample names' button.
2. Choose whether you'd like to export as a spreadsheet (xlsx) or a CSV using the checkboxes at the bottom of the interface.

## Note

Both report formats (xlsx/CSV) contain lists of taxa in single cells. In the xlsx, this is delimited by a newline (\n). In the CSV this is swapped for a semicolon ';' to avoid parsing errors.

1. Specify the output location by clicking the 'Browse' button in the 'Output File' section at the top of the interface. Fill out the 'Save As' prompt. **Please be sure add the '.xlsx' or '.csv' file extension to the filename if it not done automatically.**
2. Specify the 'Relative Abundance Threshold'. The default is 1.0%. this default parameter means that no organism < 1% relative abundance will feature in the output report.

## Important note

We ask that for the NHS RMg service evaluation, you set the Summary Report threshold cut-off to the lowest value (currently 0.1) so that all classifications are included.

1. Click 'Generate Summary'. The output will appear in the 'summary\_report' directory.

**Video tutorial - Manually selecting samples**

**Video tutorial - Providing a list of sample names**

**Output fields**

<b>Column</b>	<b>Description</b>
<b>Sample</b>	The Lab/sample ID.
<b>Experiment</b>	The specific MinNOW experiment.
<b>SampleID</b>	The specific MinNOW sample.
<b>Barcode</b>	The barcode used in the library.
<b>AnonymisedIdentifier</b>	A de-identified hospital number.
<b>CollectionDate</b>	The date on which the sample was collected.
<b>SampleClass</b>	The classification of the sample, PC, NC, NTC, specimen.
<b>SampleType</b>	The sample site, eg. BAL, SPT, NPS.
<b>Operator</b>	The name or identifier of the individual who processed or sequenced the sample.
<b>Notes</b>	Any additional notes or remarks related to the sample or its processing.
<b>Total reads XX hrs</b>	The total number of sequencing reads generated after XX hours of sequencing.
<b>Human reads XX hrs</b>	The number of reads identified as being from human DNA after XX hours of sequencing.
<b>Human reads (%) XX hrs</b>	The percentage of total reads that are identified as human DNA after XX hours.
<b>Total classified reads XX hrs</b>	The number of reads that have been classified (assigned to an organism or category) after XX hours.
<b>Sequencing N50 (bp) XX hrs</b>	The N50 statistic for the reads generated after XX hours, indicating read length distribution.
<b>Proportion &gt;Q15 quality (%) XX hrs</b>	The percentage of reads with a quality score greater than Q15 after XX hours.
<b>Median read quality (PHRED score) XX hrs</b>	The median PHRED quality score of the reads after XX hours, indicating overall data quality.
<b>Total bases (bp) XX hrs</b>	The total number of base pairs generated by the sequencing run after XX hours.
<b>Organisms (excluding viruses) XX hrs</b>	The list of organisms (excluding viruses) identified from the reads after XX hours.
<b>Organisms (excluding viruses) read counts XX hrs</b>	The read counts associated with organisms (excluding viruses) after XX hours.
<b>Organism (excluding viruses) percentage abundance XX hrs</b>	The percentage abundance of each organism (excluding viruses) in the sample after XX hours.
<b>Viral organisms XX hrs</b>	The list of viral organisms identified from the reads after XX hours.
<b>Viral read counts XX hrs</b>	The read counts associated with viral organisms after XX hours.

## 3.8 Setting up CIDR Metagenomics bioinformatics workflow - alternative deployment

---

### Note

**The tools and documentation described here and on the CIDR GitHub are not validated for diagnostic use and are for research and evaluation purposes only.**

### 3.8.1 Overview

For collaborators outside of the Network, an alternative configuration can be provided. This will bypass the GUI allowing users to provide a `sample_sheet.csv` through a CLI. Organism query will not be available to headless users as this tools is heavily reliant on GUI I/O.

### 3.8.2 Install instructions

1. Decompress CIDR\_metagenomics\_vX.X.tar.gz:

```
tar -xvzf CIDR_metagenomics_vX.X.tar.gz
```

1. Install conda/mamba.
2. Build the appropriate environment for running the CIDR metagenomics containers.

```
wget https://raw.githubusercontent.com/GSTT-CIDR/metagenomics_container/main/conda/apptainer.yml conda env create -f apptainer.yml
```

1. Allocate a directory for MinKNOW data outputs. This will be mounted to the `/data` directory in the container in a later step.

### Note

The directory structure of data for ingest must be maintained as in standard MinKNOW outputs eg. **Example for control sample**

```
[minknow_outputs_directory]/GSTT_control_sample_01/GSTT_control_sample_01/20240424_1408_X4_FAY88387_d3868a4f/fastq_paa
```

**Naming schema** `[minknow_outputs_directory]/[experiemnt]/[sample_id]/[*]/fastq_pass/barcodeXX`

### 3.8.3 Install validation

1. Navigate to the root of the `CIDR_metagenomics vX.X` directory.
2. Move `CIDR_metagenomics vX.X/GSTT_control_sample_XX` to the allocated directory for MinKNOW data outputs (from Install instructions: Step 4).

3. activate the apptainer conda environment: `conda activate apptainer`

4. Initiate the run for analysing the control dataset:

```
apptainer exec --bind .:/mnt --bind ./data:/data ./containers/cidr_metagenomics_v3.6.sif bash -c 'cd /workflow ; source /opt/conda/etc/profile.d/conda.sh'
```

1. When the workflow has completed, inspect the `CIDR_metagenomics vX.X/reports/CIDR_control_1` PDF report, it should match the CIDR validation report provided [here](#).

## Info

Variables to change in step 3

**--bind .:/mnt** - Binding the workflow root directory to the container /mnt.

**--bind ./data:/data** - binding the allocated directory for MinKNOW data outputs to /data.

**./containers/cidr\_metagenomics\_v3.6.sif** - launching the metagenomics container.

**for t in 0.5 1 2 16 24** - time-points for analysis.

**--cores 20** - number of samples to be processed simultaneously - not the same as threads.

**samples=/mnt/sample\_sheets/CIDR\_control\_1.csv** - the mounted path for the sample sheet - remember this is the relative mounted path, so `/mnt/sample_sheets` corresponds to `CIDR_metagenomics vX.X/sample_sheets` on the host machine.

## 3.8.4 Implementation

1. Build a **sample sheet** copying the structure of the example in `CIDR_metagenomics vX.X/sample_sheets`. Importantly, 'Experiment', 'SampleID' and 'Barcode' must be correct and correspond to the `[minnow_outputs_directory]/[experiment]/[sample_id]/[*]/fastq_pass/barcodeXX` scheme.

2. activate the apptainer conda environment: `conda activate apptainer`

3. Run the container, changing the flags explained in the validation step:

```
apptainer exec --bind .:/mnt --bind ./data:/data ./containers/cidr_metagenomics_v3.6.sif bash -c 'cd /workflow ; source /opt/conda/etc/profile.d/conda.sh'
```

1. PDF outputs should be found in `CIDR_metagenomics vX.X/reports/` corresponding to each LabID in the **sample sheet** loaded.

## 4. Analysis

---

### 4.1 Service evaluation report SOP

---

#### Note

Please reference use of this method in any presentation or publication as:

Alcolea-Medina, A., Alder, C., Snell, L.B. et al. Unified metagenomic method for rapid detection of microorganisms in clinical samples. Commun Med 4, 135 (2024). <https://doi.org/10.1038/s43856-024-00554-3>

#### Important

**Risk assessment for handling respiratory samples needs to be performed by each laboratory.**

**The tools and documentation described here and on the CIDR GitHub are not validated for diagnostic use and are for research and evaluation purposes only.**

## 4.2 Network validation outline

---

### Note

Please reference use of this method in any presentation or publication as:

Alcolea-Medina, A., Alder, C., Snell, L.B. et al. Unified metagenomic method for rapid detection of microorganisms in clinical samples. Commun Med 4, 135 (2024). <https://doi.org/10.1038/s43856-024-00554-3>

### Important

**Risk assessment for handling respiratory samples needs to be performed by each laboratory.**

**The tools and documentation described here and on the CIDR GitHub are not validated for diagnostic use and are for research and evaluation purposes only.**

## 5. FAQ

---

### 5.1 FAQ

---

#### 5.1.1 1) My 30 minute reports are persistently missing.

##### Presentation

During live analysis runs only, the 30-minute reports are not being produced. All other reports are coming out fine. The clock is also running an hour behind (during BST).

##### Solution

The GridION is not configured out of the factory for the UK timezone - an issue specifically during BST. this means the data ingest script timings fail and the 30 minute timepoint is missed. **Note - if other reporting timepoints are missing, this is not your solution.**

1. Type the following in to a terminal window `sudo timedatectl set-timezone Europe/London`
2. Check the time is correct. If not, [manually update the time](#) to be correct.

#### 5.1.2 2) Some timepoint reports are missing from my run.

##### Presentation

Entire samples, or timepoint reports are sporadically missing from the report folder. In many cases, reports are not generated because no microbial reads are detected - this causes the pipeline to crash. We are working on improving this error reporting and softer crashes.

There are a number of causes for reports failing. Follow the flowchart to find the correct solution.

flowchart TB A["All reports before<br> a specific time<br> point are missing. eg. <br> (0.5hr)(1hr)(2hr)(16hr)<br>(24hr)"] --> B[In the first<br> available report, is <br> the read count very low?] B -->|Yes| C[The sequencing yield<br> is low. Previous reports<br> were not generated because<br> no microbial reads were found.] B -->|No| D[See Action 1] E[A single report is<br> missing where previous<br> time points had<br> reports generated.] --> F[Action 1] G[All reports are missing] --> H[Check MinNOW barcode panel.<br> Are there reads reporting for<br> the barcode in question?] H -->|Yes| I[See Action 2] H -->|No| J[The sequencing run has<br> failed. No reads are<br> present to analyse.]

## Actions

1. Delete the folders in `../metagenomics/results` and `../metagenomics/reports` corresponding to the specific sample with issues. Rerun the the pipeline by loading the saved sample sheet. You can do this by using the 'Load Sample Sheet' function on the metagenomics launcher - find the sample sheet corresponding to the date/time. Alternatively, you can re-enter the run information in to the launcher.
2. At the moment, no report is generated if the sample consists of human reads only. Run the following command in a terminal - it opens one of the intermediate classification files which we can use to confirm a human-only run. replace `{sample_id}` with the appropriate (identical) problematic sample IDand `{timepoint}` with an appropriate timepoint eg. `0.5_hours`.

The number 9606 in the terminal output corresponds to the *Homo sapiens* taxon. If only this is present, the read is human -only. If there are other numbers, try Action 1.

## Note

Open the terminal py pressing CTRL-ALT-T

```
zcat -f /media/grid/metagenomics/results/{sample_id}/{timepoint}/centrifuge/centrifuge_raw.tsv* | awk -F'\t' '{print $3}' | sort -u
```

## 5.1.3 3) My summary report has returned a "....returned non-zero exit status 1" and is not generating spreadsheets.

---

## Presentation

The summary report tool outputs an error message and does not create a summary report.

## Actions

The tool should automatically add the correct file extension on to the filename provided by the user. If it has not, add either the '.xlsx' extension or the '.csv' extension depending on your needs.