

Jan 09, 2026

Version 2

Quality control assessment for microbial genomes: GalaxyTrakr MicroRunQC workflow V.2

 Version 1 is forked from [Quality control assessment for microbial genomes: GalaxyTrakr MicroRunQC workflow](#)



DOI

<https://dx.doi.org/10.17504/protocols.io.261ge138dv47/v2>

Candace Bias¹, Ruth Timme¹, Yesha Shrestha², Tina Pfefer³, Paul Morin⁴, Maria Balkey³, Errol Strain³

¹US Food and Drug Administration; ²Center for Veterinary Medicine, US Food and Drug Administration;

³Center for Food Safety and Applied Nutrition, U.S. Food and Drug Administration, College Park, Maryland, USA;

⁴U.S. Food and Drug Administration, Jamaica, New York, USA

GenomeTrakr

Tech. support email: genomeTrakr@fda.hhs.gov



Maria Balkey

US Food and Drug Administration

Create & collaborate more with a free account

Edit and publish protocols, collaborate in communities, share insights through comments, and track progress with run records.

[Create free account](#)

OPEN  ACCESS



DOI: <https://dx.doi.org/10.17504/protocols.io.261ge138dv47/v2>

Protocol Citation: Candace Bias, Ruth Timme, Yesha Shrestha, Tina Pfefer, Paul Morin, Maria Balkey, Errol Strain 2026. Quality control assessment for microbial genomes: GalaxyTrakr MicroRunQC workflow. **protocols.io**
<https://dx.doi.org/10.17504/protocols.io.261ge138dv47/v2>Version created by **Maria Balkey**

License: This is an open access protocol distributed under the terms of the **Creative Commons Attribution License**, which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited

Protocol status: Working

We use this protocol and it's working

Created: January 09, 2026

Last Modified: January 09, 2026

Protocol Integer ID: 238313

Keywords: WGS, Quality Control, GalaxyTrakr, GenomeTrakr, microbial pathogen survielliance, galaxytrakr microrunqc workflow, galaxytrakr microrunqc workflow purpose, quick access to genometrakr sequence quality threshold, wgs sequence quality for bacterial pathogen, quality control assessment for microbial genome, genometrakr sequence quality threshold, check against genometrakr qc, microrunqc workflow, microrunqc, genometrakr qc, microbial genome, checking wgs sequence quality, quality assessments for raw read, genometrakr, microbial pathogen, galaxytrakr, most microbial pathogen, bacterial pathogen, sequence type for each isolate, sequence type definition file, de novo assembly, available in sequence type definition file, end fastq file, pathogen, galaxytrakr account, raw read, cronobacter threshold, account in galaxytrakr, custom galaxy instance, mlst method, quality control assessment, assembly qc, added enterobacter qc, additional mlst data field, galaxytrakr upgrade, nextseq, entire miseq

Disclaimer

Please note that this protocol is public domain, which supersedes the CC-BY license default used by protocols.io.

Abstract

PURPOSE: Step-by-step instructions for checking WGS sequence quality for bacterial pathogens. The MicroRunQC workflow, implemented in a custom Galaxy instance, will produce quality assessments for raw reads (Illumina paired-end fastq files) and draft de novo assemblies, along with reporting the sequence type for each isolate. This workflow will work on most microbial pathogens, so we advise laboratories to upload their entire MiSeq/NextSeq run through this workflow.

SCOPE: This protocol covers the following tasks:

1. Quick access to GenomeTrakr sequence quality thresholds by organism
2. Create a GalaxyTrakr account
3. Set up an account in GalaxyTrakr
4. Create a new history/workspace
5. Upload data
6. Execute the MicroRunQC workflow
7. Interpret the results - check against GenomeTrakr QC thresholds

Version updates:

V7: Edits to incorporate GalaxyTrakr upgrades and new interface.

V6: Minor edits, including section reorganization and addition of clarifying notes

V5: New column in the output table to capture additional mlst data fields when available in Sequence Type definition files (not available for all species)

V4: MicroRunQC updated to V1.1 Includes updates to skeza and mlst methods, as well as adjusted assembly QC thresholds for E.coli. Added *Enterobacter* QC thresholds to threshold table.

V3: updated with *Cronobacter* thresholds

Troubleshooting

Quick Access to QC Benchmarks

- 1 This protocol will walk the user through various aspects of the quality assessment of bacterial genome sequences, from setting up a GalaxyTrakr account to the quality control (QC) benchmarks GenomeTrakr uses for its sequencing efforts. For quick access, GenomeTrakr QC benchmarks are included in the table below.

These are also relevant for NARMS and VetLIRN contributors.

*MicroRunQC users should follow QC threshold guidelines established by their respective surveillance coordinating body(s).

	A	B	C	D	E	F	G	H	I	J
Quality metric		<i>Salmonella</i>	<i>Listeria</i>	<i>E. coli</i>	<i>Shigella</i>	<i>Campylobacter</i>	<i>Vibrio para.</i>	<i>Cronobacter</i>	<i>Enterococcus faecium</i>	<i>Enterococcus faecalis</i>
Average read quality Q score for R1 and R2		>=30	>=30	>=30	>=30	>=30	>=30	>=30	>=30	>=30
Average coverage		>=30X	>=20X	>=40X	>=40X	>=20X	>=40X	>=20X	>=50X	>=40X
<i>De novo</i> assembly: Seq. length (Mbp)		~4.3-5.2	~2.7-3.2	~4.5-5.9	~4.0-5.0	~1.5-1.9	~4.8-5.5	~4-5	~2.5-3.5	~2.5-3.25
<i>De novo</i> assembly: no. contigs		<=300	<=300	<=400	<=550	<=300	<=300	<=500	<=350	<=200

Account set up

- 2 1. Create a GalaxyTrakr account here: <https://account.galaxytrakr.org/Account/Register>

Note

This is a more detailed form than what is available by clicking "Register here" on the genometrakr.org main page. Please use the form linked here when creating a new account.

User Registration Form

Location Add New Location

First Name

Last Name

Email
Email will be used for automated messages to include registration information!

Primary Phone
If possible please use a mobile number than can accept text messages, only used for support

Title

Requirements

2.1 Log into your GalaxyTrakr account: <https://galaxytrakr.org>

Create a new history

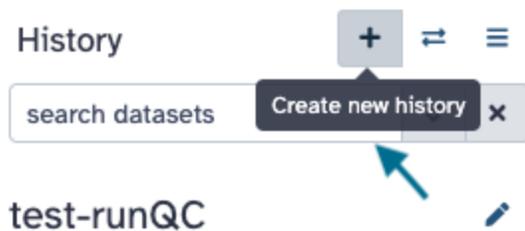
3 Create a new history.

We recommend creating a new history for each new MiSeq Run and including the flow-cell ID and date in the history name.

Save your MicroRunQC output here and any other relevant analyses, like serotyping, or AMR detection.

After all the analysis output from this run is saved to your internal data network or computer, older histories should be purged/deleted so as not to occupy the limited storage space in your account. In some cases it may be useful to save, for a limited time, multiple histories or to run analyses concurrently in multiple histories. In these cases you need to pay attention to your % usage bar (shows % used of allocated storage space) in the upper right corner of the GalaxyTrakr page. If you need additional space you can contact galaxytrakrsupport@fda.hhs.gov and request additional storage.

3.1 Click on the + icon in the upper right History panel



3.2 Name your new History by clicking on the "Unnamed history" text, type in desired name, and hit Enter. We recommend including the run cell ID and the date the run was started.

History





Unnamed history



Add Tags



B

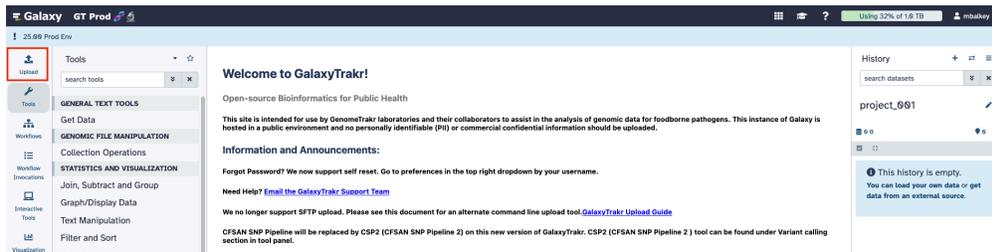


This history is empty.

You can load your own data or get data from an external source.

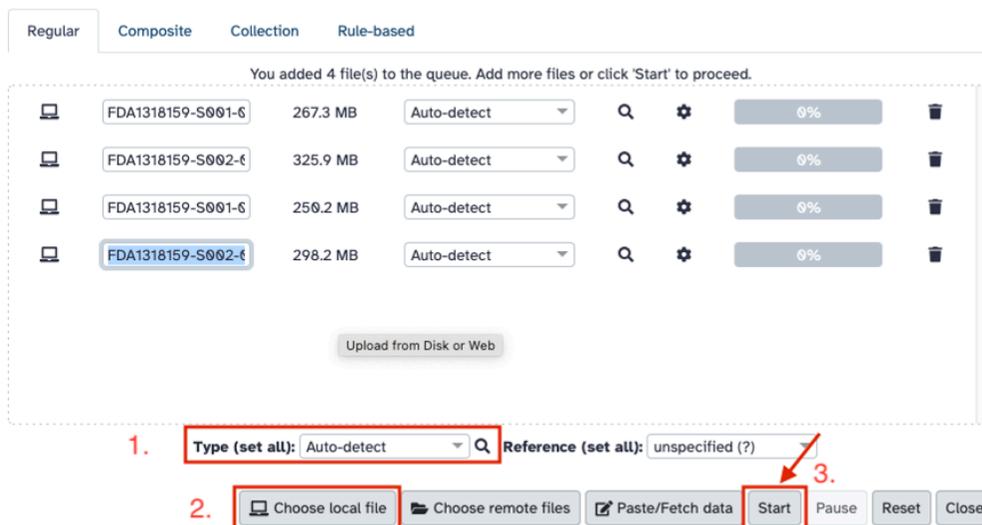
Upload data

- 4 **This section will describe the process for uploading raw fastq files into your active History panel.** After the files have been uploaded they will stay in your account until they are deleted.
 - 4.1 Click on the Upload icon in the GalaxyTrakr menu to start an upload process.



- 4.2 Select "**Type (set all):auto-detect.**" Click "**Choose local files**" button and navigate to the desired fastq files, then click "**Start**" to upload files. These files should be paired (two per sample/isolate).

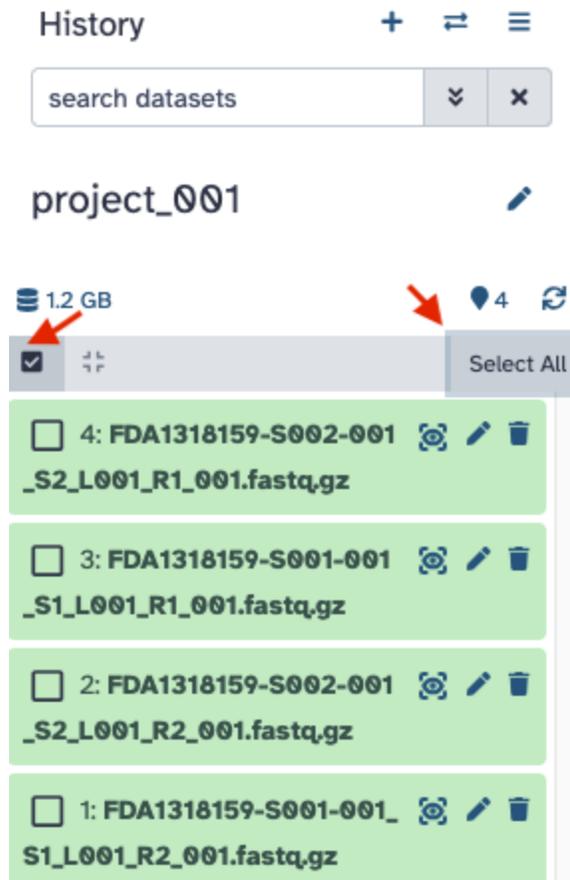
Upload from Disk or Web to **project_001**



- As the file uploads complete, each row will turn green. Samples in yellow are still in process.
- 4.3 You have just upload a set of forward and reverse reads. For further analysis, these files need to be paired properly so the platform knows which R1 and R2 files go together.

GalaxyTrakr does this by creating a **List of Dataset Pairs**.

Within your newly created History panel, click the check mark box, then select all the files you just uploaded.



Screenshot of History panel showing recently uploaded files. Note the way the files are named, using R1 and R2 to identify the paired reads. This will be important in the next step. Some naming conventions can be slightly different.

- 4.4 Click "Select All" and choose "Advance Build List", select "List of Paired Datasets"

4.5 A new window will open to help you pair the fastq files properly. By default `_R1` and `_R2` are the selected options to pair fastq files. Note how your paired reads are named.

Click on the filter icon if "`_R1`," "`_R2`" have to be replaced by other suffix.

Click **Next**, double check the name of the datasets. If the names are accurate, name the List of Paired Datasets and click **Build**.

What are you building? ✓ Auto Pairing ✓ 3 Builder

Assemble, label, and sort your list of pairs.

This interface allows you to build a new Galaxy list of pairs. List of pairs are an ordered list of individual datasets paired together in their own paired collection (often forward and reverse reads).

Auto-matched 2 pair(s) of datasets from target datasets. If this isn't correct, [configure auto-pairing](#).

Dataset(s)	List Identifier	Discard	Status
FORWARD 4: FDA1318159-S002-001_S2_L0C REVERSE 2: FDA1318159-S002-001_S2_L001	FDA1318159-S002-001_S2_L001_001	✗	✓
FORWARD 3: FDA1318159-S001-001_S1_L00 REVERSE 1: FDA1318159-S001-001_S1_L001	FDA1318159-S001-001_S1_L001_001	✗	✓

Remove file extensions? Hide original elements

Name:

Back Build

History

search datasets

project_001

1.2 GB

All 4 selected

- 4: [FDA1318159-S002-001_S2_L001_R1_001.fastq.gz](#)
- 3: [FDA1318159-S001-001_S1_L001_R1_001.fastq.gz](#)
- 2: [FDA1318159-S002-001_S2_L001_R2_001.fastq.gz](#)
- 1: [FDA1318159-S001-001_S1_L001_R2_001.fastq.gz](#)

4.6 This List of Paired Datasets will be available for analysis in your history panel. You can run multiple analyses on the same dataset rather than upload the same sequence data to a new history to perform additional analyses. This will help you use your allocated storage space efficiently.

History

search datasets

project_001

1.2 GB

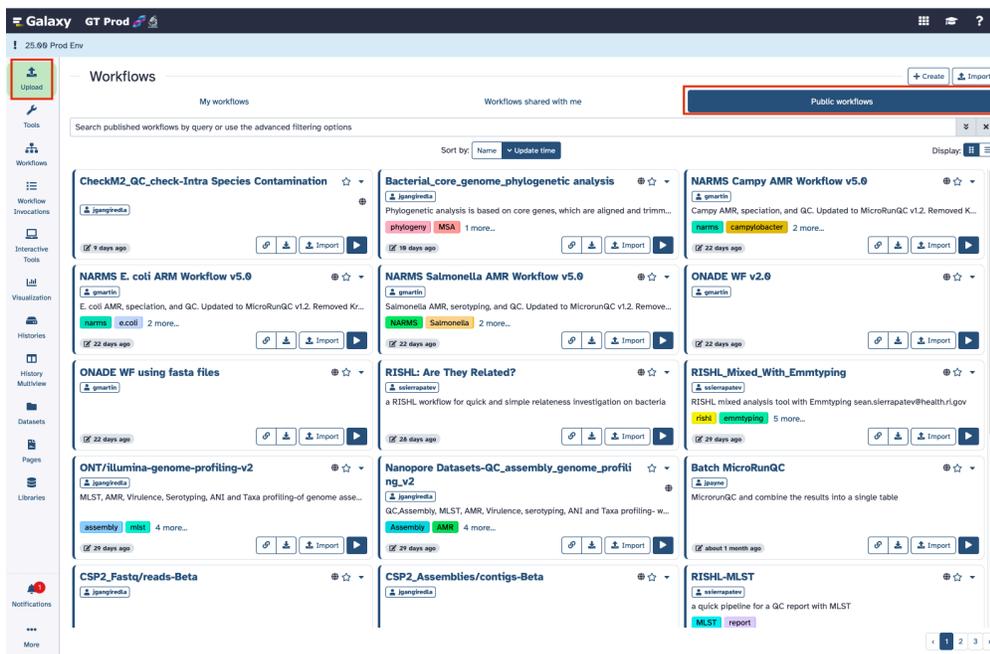
1 8

9: project_001

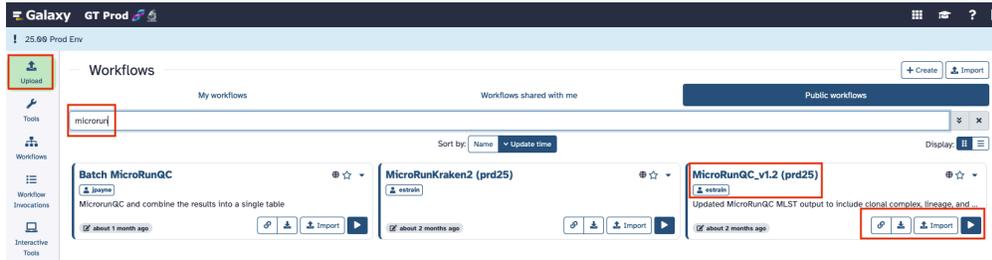
a list with 2 fastqsanger.gz pairs

Run the MicroRunQC workflow

- 5 Add the MicroRunQC workflow to your own **"Workflows"** panel. You only have to do this step once for each new workflow you need.
- 5.1 Navigate to the **"Workflow"** tab on the main menu, click **"Public workflows"**, and search for "MicroRunQC_v1.2."

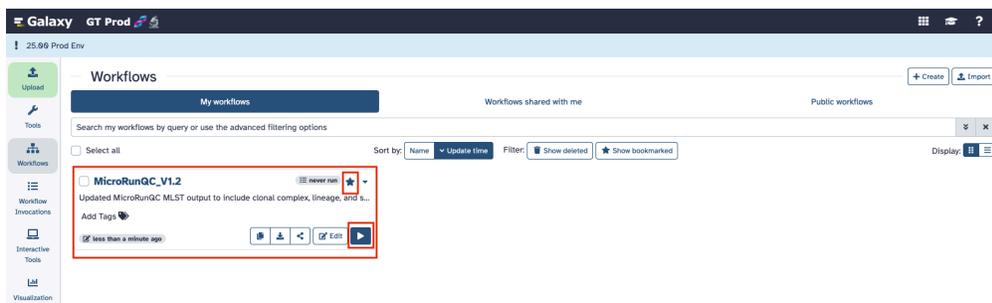


Click **"Import"** to select MicroRunQC.

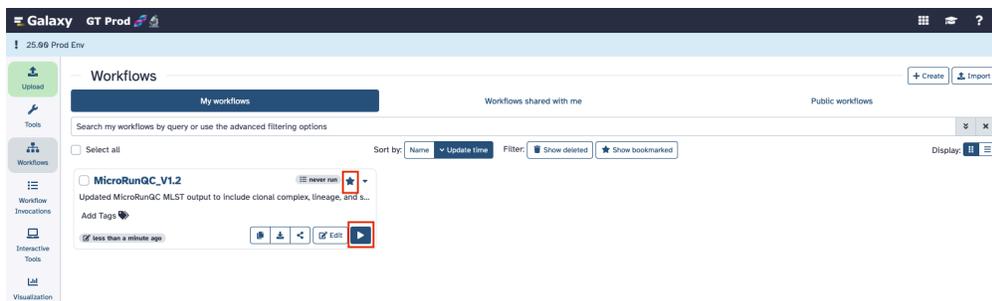


5.2 To see the new imported workflow, click the **"Workflow"** tab on the main menu, click "My Workflows."

Click the star to bookmark this workflow.

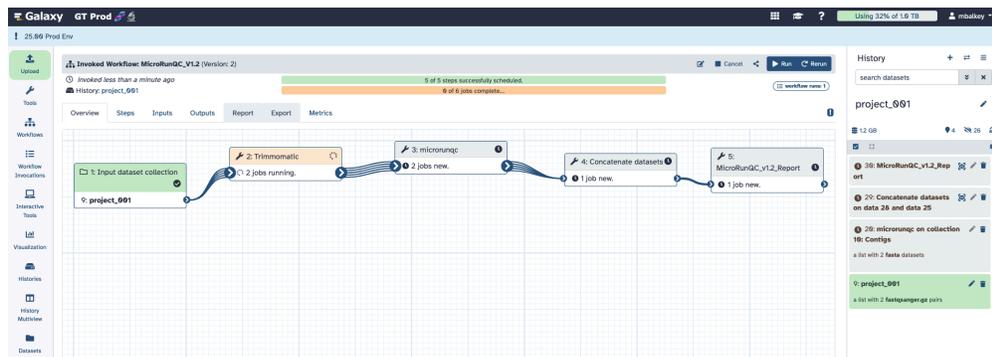


5.3 Click the play icon to run the **MicroRunQC_v1.2** and select the dataset collection that was created earlier.



5.4

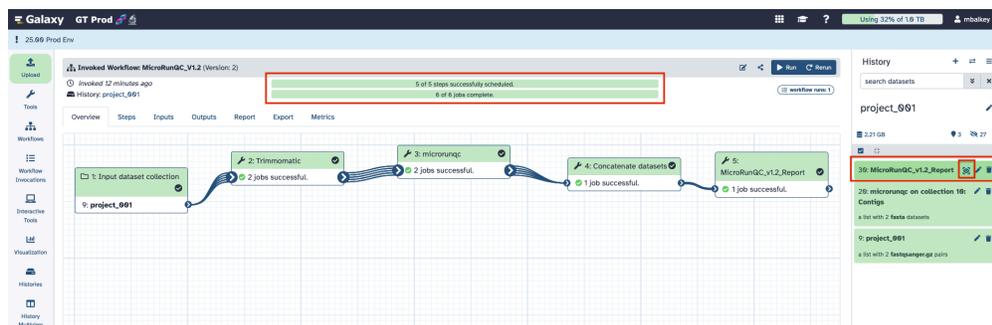
Click **Run Workflow**. This can take some time depending on the number of samples you are analyzing. If you choose to you can log out of GalaxyTrakr and log back in at a later time to see if the job is completed.



5.5

Upon completion of the pipeline all tiles in the History pane will be green and each of the steps in the pipeline will show "completed".

In the **"MicroRunQC_v1.2_Report"** tile, click on the **"Eye"** icon to view the output table in the GalaxyTrakr window.



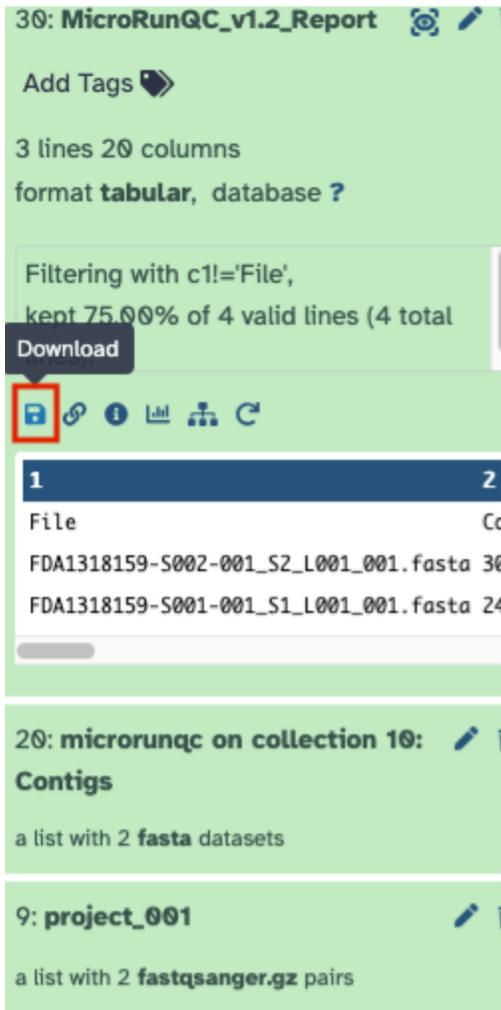
Interpret the results

6 Download and interpret the results:

- 6.1 Click **MicroRunQC_v1.2_Report** and then the floppy disc icon. The tabular file can be opened in a text reader or converted to a format (.txt) that can be opened in Excel.

The screenshot shows the Galaxy web interface with a table titled "36: MicroRunQC_v1.2_Report". The table has 15 columns and 2 rows of data. The columns are labeled as follows: Column 1 (File), Column 2 (Contigs), Column 3 (Length), Column 4 (EstCov), Column 5 (NS9), Column 6 (MedianInsert), Column 7 (MeanLength_R1), Column 8 (MeanLength_R2), Column 9 (MeanQ_R1), Column 10 (MeanQ_R2), Column 11 (Scheme), Column 12 (ST), Column 13, Column 14, Column 15, Column 16, Column 17, Column 18, and Column 19. The data rows are as follows:

Column 1	Column 2	Column 3	Column 4	Column 5	Column 6	Column 7	Column 8	Column 9	Column 10	Column 11	Column 12	Column 13	Column 14	Column 15	Column 16	Column 17	Column 18	Column 19
File	Contigs	Length	EstCov	NS9	MedianInsert	MeanLength_R1	MeanLength_R2	MeanQ_R1	MeanQ_R2	Scheme	ST							
FDA1318159-5682-661_L52_L661_661.fasta	39	3892474	295.1	325318	152	121.8	122.8	36.6	35.6	listfa_2	161	CC+CC181Llineage=11	abc2(7)	bg(A15)	cat(5)	dup(8)	del(6)	idh(14)



6.2 The MicroRunQC output file includes the following columns:

A	B	C
Parameter	Input	Description
Contigs	Assembly	Number of contigs in the de-novo SKESA assembly. Contigs smaller than 200 base-pairs (bp) are not counted.
Length	Assembly	Total length of all contigs > 200bp. This should approximate the size of the genome for the target organism.
EstCov	Assembly	Mean coverage for contigs in the SKESA assembly.



A	B	C
N50	Assembly	Sequence length of the shortest contig at 50% of the total genome length
MedianInsert	Read	Distance between forward and reverse reads. Calculated by mapping reads to SKESA assembly using bwa.
MeanLength_R1	Read	Mean length of forward read
MeanLength_R2	Read	Mean length of reverse read
MeanQ_R1	Read	Mean Q-score of forward read
MeanQ_R2	Read	Mean Q-score of reverse read
Scheme	Assembly	PubMLST scheme name (output from mlst application that scans contig files against traditional PubMLST typing schemes).
ST	Assembly	Sequence Type
MLST extra	Assembly	e.g. <i>Listeria</i> clonal complex info
Loci	Assembly	gene (allele number) – for example aroC(118)

MicroRunQC output table headers. This table lists the summary metrics for sequence quality, number of contigs, and estimated genome size, along with other common metrics for reads (Median Insert Size and Mean Length) and assemblies (N50). Additionally, if the Multi-Locus Sequence Type (MLST) for the isolate is available from pubmlst, the workflow also reports Sequence Type (ST) and the associated alleles.

***MLST extra:** Additional data fields reported when available in Sequence Type definition files (not available for all species)

1. clonal_complex – sequences grouped by similarity to central allelic profile (e.g., *Campylobacter* ST-21 complex)
2. CC – clonal_complex – Abbreviation used for organism like *Listeria*, ST profiles are maintained by different groups
3. Lineage – *Listeria monocytogenes* lineage (I,II,III, and IV), *Listeria* species also reported here (e.g. *L.innocua*)
4. species – e.g., *Vibrio alginolyticus*

**This output should be saved either to your LIMS or to a spreadsheet linked to the sequencing run and samples.

6.3 Example output for 1 *Salmonella* and 5 *Listeria* isolates.

A	B
Srain ID	Lab Confirmation
FDA1216271-C001-001	Listeria mono
FDA817806-S073-001	Listeria mono
FDA746634	Listeria mono
FDA1213377-C001-002	Listeria grayi
FDA933376-S060-005	Listeria innocua
FDA1213835-C001-001	Salmonella

Lab confirmed IDs for 6 isolates

A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T
File	Contigs	Length	EstCov	N50	MedianInsert	MeanLength $\bar{R}1$	MeanLength $\bar{R}2$	MeanQ $\bar{R}1$	MeanQ $\bar{R}2$	Scheme	ST	MLSTextra							
FDA1216271-C001-001	16	2911949	36.7	476210	321	148.4	148.4	36.4	34.6	listeria_2	5	CC=CC5;Lineage=1	abcZ(2)	bgIA(1)	cat(1)	dapE(3)	dat(3)	ldh(1)	lhkA(7)



	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T
	FD A 8 1 7 8 0 6 - S 0 7 3 - 0 0 1	2 0	3 0 6 8 3 5 4	1 7 9 . 6	5 2 5 4 3 8	3 2 9	2 3 4 . 7	2 3 5 . 2	3 6 . 7	3 1 . 9	li st e ri a ₂	3 2 1	C C = C C C 3 2 1 , L i n e a g e = 1	a b c Z (5)	b g l A (6)	c a t (8)	d a p E (6 2)	d a t (6)	l d h (7)	l h k A (3 4)
	FD A 7 4 6 6 3 4	3 0	3 0 5 2 8 8 8	4 1 . 4	2 9 3 9 4 7	3 2 0	1 4 8 . 4	1 4 8 . 4	3 6 . 5	3 6	li st e ri a ₂	-		a b c Z (2)	b g l A (1)	c a t (1)	d a p E (3)	d a t (3)	l d h (1)	l h k A (~ 7)
	FD A 1 2 1 3 3 7 7 - C 0 0 1 - 0 0 2	2 0	2 6 7 2 1 8 0	1 5 5 . 1	4 7 3 1 8 1	2 7 0	1 4 7 . 3	1 4 7 . 3	3 7 . 2	3 6 . 1	-	-								
	FD A 9 3 3 3 7 6 - S	9	2 8 8 1 8 6 9	2 1 3	1 4 9 8 7 9 0	3 0 3	2 3 2 . 1	2 3 2 . 2	3 7	3 6 . 2	li st e ri a ₂	1 4 8 9	C C = C C 1 4 8 9 , L i n	a b c Z (2 5 0)	b g l A (2 1)	c a t (8 3)	d a p E (2 9 8)	d a t (2 0)	l d h (4 5 8)	l h k A (2 1 6)

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T
	0 6 0 - 0 0 5												e a g g e = L. i n n o c u a							
	F D A 1 2 1 3 8 3 5 - C 0 0 1 - 0 0 1	3 7	4 8 3 2 3 6 5	3 4 .4	2 9 4 9 3 6	3 5 4	1 4 9	1 4 9	3 6. 6	3 5. 7	s e n t e r i c a - a c h t m a n - 2	2 1 4		a r o C (1 4)	d n a N (7 2)	h e m D (2 1)	h i s D (1 2)	p u r E (6)	s u c A (1 9)	t h r A (1 5)

MicroRunQC example report showing mlst ST results for different *Listeria* species.

The mlst *Listeria* database includes multiple species, including *Listeria monocytogenes* and *L. innocua*. When available, the *Listeria* clonal complex (CC) or *L. monocytogenes* lineage is listed alongside the ST.

6.4 For quality control threshold guidelines for the GenomeTrakr surveillance network,

 [go to step #1](#) These are also relevant for NARMS and VetLIRN contributors.

*MicroRunQC users should follow QC threshold guidelines established by their respective surveillance coordinating body(s).