



Apr 06, 2022

Wastewater QC workflow in GalaxyTrakr (SSQuAWK3) V.5

Jasmine Amirzadegan¹, Tunc Kayikcioglu¹, hugh.rand¹, Ruth Timme², Maria Balkey¹

¹Center for Food Safety and Applied Nutrition, U.S. Food and Drug Administration, College Park, Maryland, USA;

²US Food and Drug Administration



dx.doi.org/10.17504/protocols.io.kxygzk5dv8j/v5

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

Jasmine Amirzadegan

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

PURPOSE:

Step-by-step instructions for checking sequence quality for SARS-CoV-2 wastewater samples using **SSQuAWK3: SARS - CoV - 2 Sequence Quality Assurance Workflow and Kontraption, version 3**. The SSQuAWK3 workflow, implemented in CFSAN's custom Galaxy instance (GalaxyTrakr) will produce quality assessments for raw reads (Illumina MiSeq paired-end fastq files).

SCOPE: This protocol covers the following tasks:

1. Set up an account in GalaxyTrakr
2. Create a new history
3. Upload data and reference files
4. Execute the SSQuAWK3 workflow
5. Interpret the results

Protocol and SSQuAWK workflow version history:

Protocol V1, SSQuAWK version 1: Basic protocol steps with screenshots

Protocol V2, SSQuAWK version 1: Addition of a detailed 12 minute video tutorial

Protocol V3, SSQuAWK version 2: Addition of 5 new genome mapping metrics

Protocol V4 SSQuAWK version 3: Metrics now reported with fewer softwares, fewer underlying GalaxyTrakr jobs, and about 50% fewer underlying GalaxyTrakr steps. Cleaner output table formats now include QC placeholder columns for SRA metadata template.

Protocol V5 SSQuAWK version 3: Previous protocol version had broken links for FASTA and BED files, this version fixes the links.

DOI

dx.doi.org/10.17504/protocols.io.kxygzk5dv8j/v5

<https://galaxytrakr.org>

Jasmine Amirzadegan, Tunc Kayikcioglu, hugh.rand , Ruth Timme, Maria Balkey 2022. Wastewater QC workflow in GalaxyTrakr (SSQuAWK3). **protocols.io**
<https://dx.doi.org/10.17504/protocols.io.kxygzk5dv8j/v5>
Jasmine Amirzadegan

WGS, Quality Control, GalaxyTrakr, GenomeTrakr, microbial pathogen surveillance

protocol ,

Apr 06, 2022

Apr 06, 2022

60411

:

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

Account set up

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

User Registration Form

Location	California Department of Public Health - Food and Drug Laboratory Branch
	Add New Location
First Name	<input type="text"/> <small>Enter First Name. Do not use characters: / ?<=>~*~"~'~
~</small> </small>
Last Name	<input type="text"/> <small>Enter First Name. Do not use characters: / ?<=>~*~"~'~
~</small> </small>
Email	<input type="text"/> <small>Email will be used for automated messages to include registration information!</small>
Primary Phone	<input type="text"/> <small>Please enter number with country codes, without dashes, for example +17085520189</small>
	<small>If possible please use a mobile number then we can accept text messages, only used for support</small>
Title	<input type="text"/>
Requirements	<small>Please provide intended use of Galaxy and Analytic tools. List specific tools you would like to see deployed in Galaxy.</small>
	<input type="button" value="Register"/>

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

Galaxy / GalaxyTrakr 1905

Analysis Data Workflow Visualize Shared Data Help Login

Welcome to Galaxy, please log in

Username or Email Address

Password

Forgot password? Click here to reset your password.

Log In

Don't have an account? Registration for this Galaxy instance is disabled. Please contact an administrator for assistance.

Welcome to GalaxyTrakr: open-source bioinformatics for public health.

This site is intended for use by GenomeTrakr laboratories and their collaborators to assist in the analysis of genomic data for foodborne pathogens.

This instance of Galaxy is hosted in a public environment and no personally identifiable (PII) or commercial confidential information should be uploaded.

--||--Information and Announcements--||--

Please re-improve the skemastran workflow that was updated a few days ago. Previous versions are no longer working and are causing errors when running. Thank you.

Access CFSPAN SNP Pipeline workflows in the shared workflows screen.

Post in the official Galaxy GenomeTrakr board on the Redmine site. Click here

Click here to access the GalaxyTrakr User Guide

Forgot Password? Email GalaxyTrakr Support Team

Create a new history

2 Create a new history.

We recommend creating a new history for each new MiSeq sequence set with details and date in the history name.

Save your SSQuAWK output here with any other relevant analyses.

After all the analysis output from this run is saved to your internal data network or computer, older history's 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.

2.1 Create a new history with the "+" symbol in the upper right hand corner. Name your history and press "enter" on your keyboard to save the name.

[illegible]

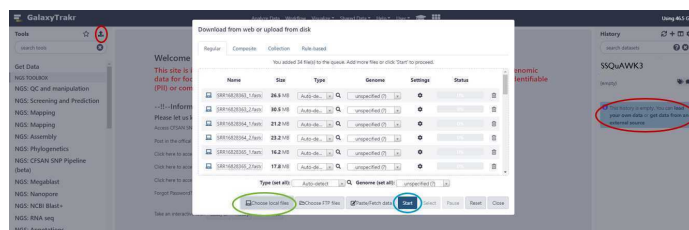
Upload sequence data

3 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.

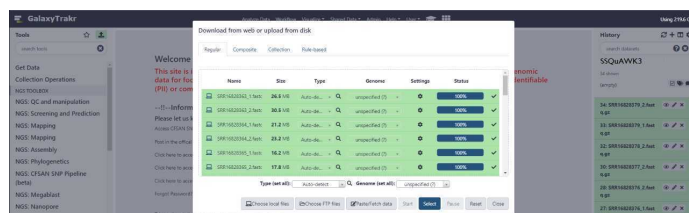
3.1 Upload sequence data to your history, using either of the two options circled in red below.

A window will appear in the middle of your screen. This is where you select your files using the "Choose local files" button at the bottom of the window. The "Choose local files" button is circled in green. These fastq files should be paired (two per sample).

After you've selected your files, press "Start" to initiate your data upload to GalaxyTrakr. The "Start" button is circled in blue.



3.2 As the file uploads complete, each row will turn green. If samples are shown with yellow background, then are still uploading.



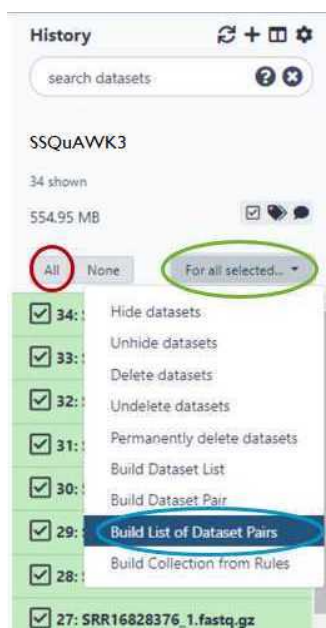
3.3 You have just upload a set of forward and reverse reads. For further analysis these files need to be paired properly so the platform knows which R1 and R2 files go with each sample. GalaxyTrakr does this by creating a **List of Dataset Pairs**.

Within your newly created History panel, click the "check box," then select all the files you just uploaded by clicking "All" or by individually selecting the ones you want to pair.



- 3.4 Check all the files belonging to a pair. In this example, all the files belong to a pair, so I will use the "All" button (circled in red).

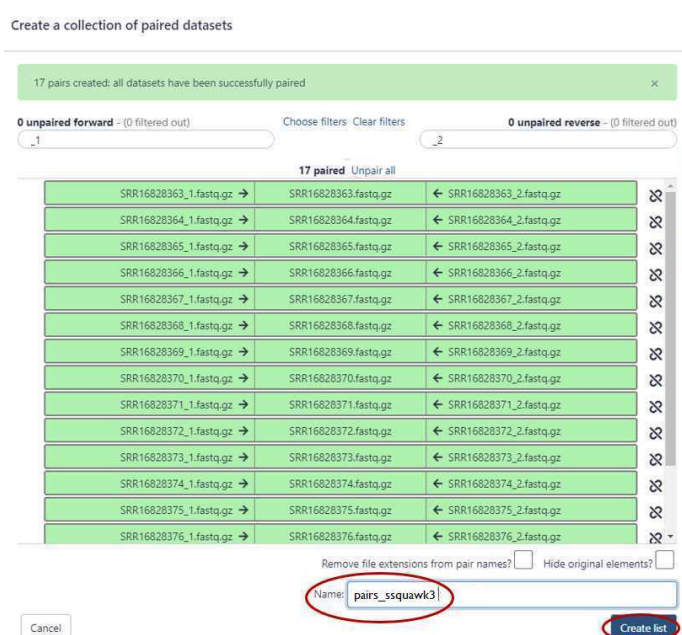
Then, use the "For all selected..." dropdown (circled in green), and click on "Build List of Dataset Pairs" (circled in blue).



- 3.5 GalaxyTrakr will automatically pair the files, but it's good to double check.

Paired reads will pair in the middle column and turn green.

If everything looks good, then choose a name for your pairs (circled red) and "Create List" (also circled red).



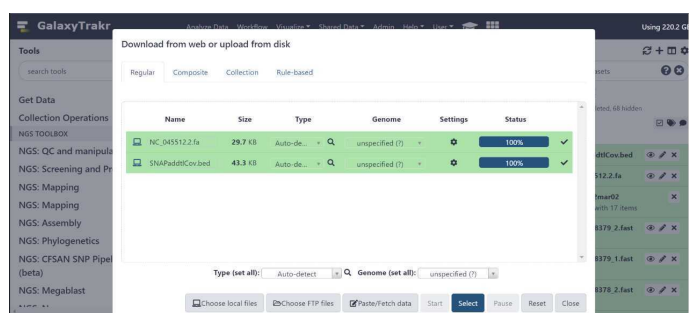
Alternatively, instead of auto-pairing you can click "choose filters" and select the appropriate filter for the pairing:



- 3.6 This paired dataset will now be available for analysis in your history panel. You can run multiple analyses on the same dataset in a history 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.



- 4 To the existing history, also upload (1) the **provided reference.fasta file** and (1) a primer.bed file.



- 4.1 SSQuAWK2 is only compatible with the 22903 nt reference genome file obtained from NCBI 'NC_045512.2'. It is provided here for your convenience:

[NC_045512.2.fa](#)

- 4.2 The primer.bed file should correspond to the SARS - CoV - 2 enrichment primer panel kit used.

QIAseq Direct kit: [QIAseqDIRECT.bed](#)

SNAP standard kit: [SNAPStd.bed](#)

SNAP additional coverage kit: [SNAPaddtlCov.bed](#)

NEB VarSkip Short (version 1) kit: [varskipShort.bed](#)

NEB VarSkip Short (version 2) kit: [VSSv2.primer.bed](#)

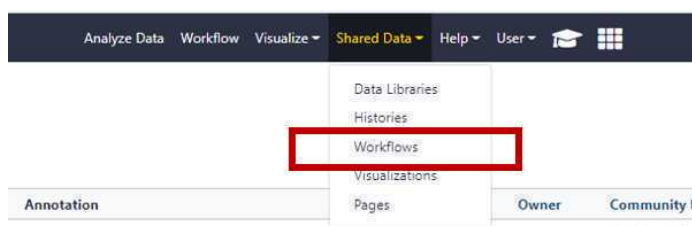
ARTIC v4 primer schemes: [ARTICv4.bed](#)

Run the SSQuAWK workflow

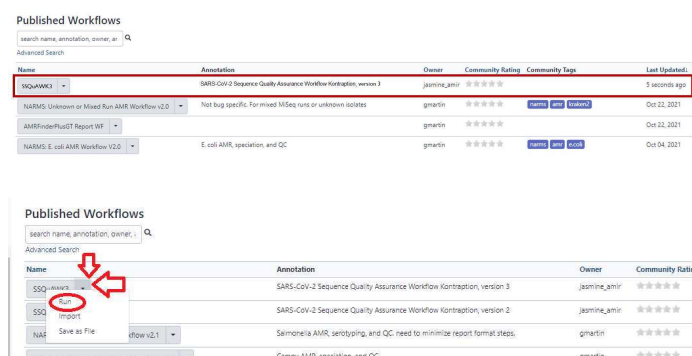
5 Access the SSQuAWK3* workflow with the "workflows" panel.

*SSQuAWK3: SARS - CoV - 2 Sequence Quality Assurance Workflow Kontrapion, version 3

- 5.1 Navigate to the "Shared Data" drop down and choose workflows



Then, from the SSQuAWK3 drop down menu, select "Run".



- 5.2 Select the paired list you created earlier by selecting the folder icon (boxed in red), and then the list

of pairs (boxed in green).

Boxed in gold: Select the reference fasta file from your history.

Boxed in blue: Select the bed file from your history.

Click Run Workflow (boxed in purple).



Running the workflow can take some time depending on the number of samples you are analyzing. Once GalaxyTrkr adds the workflow invocation to the queue, you can choose to log out of GalaxyTrkr and log back in at a later time to see if the job is completed.

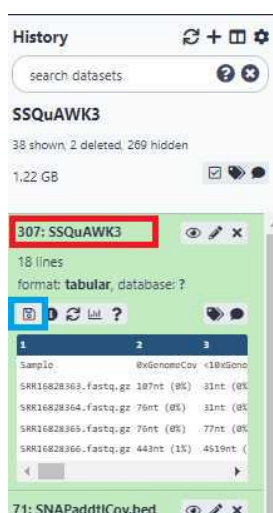


- 5.3 Upon completion of the pipeline, the output file for **SSQuAWK3** will be green. Click on the “Eye” icon to view in GalaxyTrkr window.

Interpret the results

6 Download and interpret the results:

- 6.1 Click the output file text for “**SSQuAWK3**” (boxed in red) and then the floppy disc save icon (boxed in blue). The tabular file can be opened in a text reader or converted to a format that can be opened in Excel.

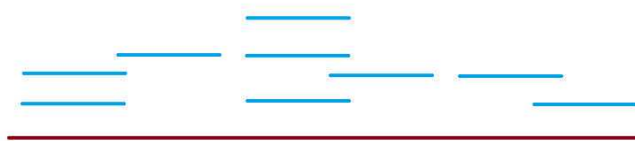


6.2 The SSQuAWK3 output file includes the following metrics:

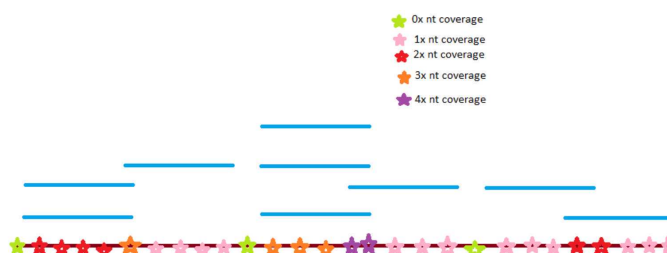
A	B	C
Parameter	Input	Description
Sample	List of Pairs	Sample name from list of pairs
0xGenomeCov	Bowtie2, samtools, ivar_trim	Percentage of nucleotides that do not cover the genome at all (zero times)
<10xGenomeCov	Bowtie2, samtools, ivar_trim	Percentage of nucleotides that barely cover the genome (less than 10 times)
nReads	Bowtie2	Total number of reads
avgLen	Bowtie2, samtools	Average read length
avgLenPassFilt	Bowtie2, samtools, ivar_trim	Average read length after iVar_trim filtering*
avgQual	Bowtie2, samtools	Average read quality
avgQualPassFilt	Bowtie2, samtools, ivar_trim	Average read length after iVar_trim filtering*
avgCovPassQual	Bowtie2, samtools, ivar_trim	Average number and percentage of nts from sequence reads that map to the genome
readsAlign	Bowtie2, samtools	Number and percentage of reads that aligned to the reference sequence.
readsAlignPassFilt	Bowtie2, samtools, ivar_trim	Number and percentage of reads that aligned to the reference sequence after iVar_trim filtering*.
humanReads	Kraken2	Number and percentage of reads classified as <i>Homo sapiens</i>
SARS-CoV-2Reads	Kraken2	Number and percentage of reads classified as SARS - CoV - 2
syntheticSeqsReads	Kraken2	Number and percentage of reads classified as non - biological sequences
quality_control_method_name	SSQuAWK	Name of the method or pipeline used to evaluate sequence quality
quality_control_method_version	3.0	Version number of the quality control pipeline or method used
quality_control_determination		Result of the quality control assessment. Blank if pass/fail thresholds have not been established or "sequence flagged for potential quality control issues" if relevant.
quality_control_issues		More information for sequences that have a QC flag issue

* The iVar_trim filter parameters: minReadLen = 30, minQual_slidingWindow = 20, and slidingWindow = 4 nt.

6.3 What is nucleotide coverage?! Let's look at 2 simple pictures



In the figure above, let the burgundy line represent the entire reference genome.
The blue lines are the reads, as sequenced nucleotides.



In the figure above, each star, drawn on the burgundy line (reference genome) is a **nucleotide position**.

There are 28 stars, so we will say our genome is 28 nucleotides long.

We can use coverage to determine the quality of our sequences (blue lines).

The lime green stars along the genome represent 0X coverage, because we did not sequence any reads with **nucleotides positions covering that reference nucleotide position**. There are no blue lines that we sequenced there!

There are 3 nucleotide positions with 0x coverage. The total genome is 28 nucleotides long.

$$\begin{aligned}\text{percent_nt0Xcov} &= (\text{nucleotidePositions0Xcov} / \text{genomeLength}) * 100 \\ \text{percent_nt0Xcov} &= (3 / 28) * 100 \\ \text{percent_nt0xcov} &= 10.71\%\end{aligned}$$

In most ideal scenarios, higher coverage indicates better sequence quality.

For example, 100x coverage is better than 10x coverage.

Since we want **higher coverage**, percent_nt0Xcov and percent_ntLess10Xcov are ideally **lower percentages**.

0x coverage and 10x coverage indicate "no coverage" and "poor coverage", respectively.

Generally, we expect avgReadCov in 100's or 1000's*

If **percent_nt0Xcov** is a higher percentage, say 50%*, that means half of the genome was not covered by our sequences. The quality is not ideal.

** These values are not official threshold and only used for illustrative purposes.*

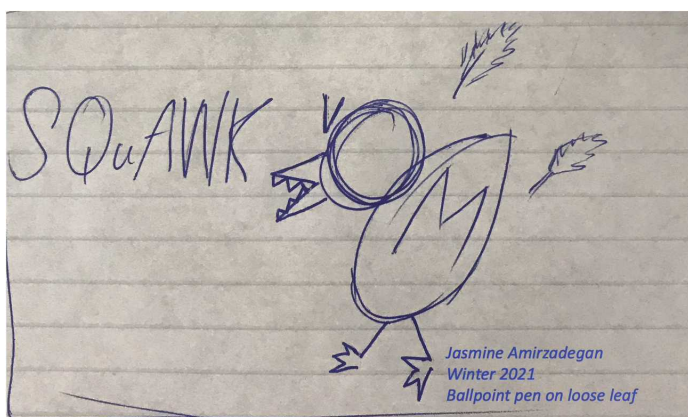
Threshold guidance is 'in progress', and planned to be announced after further analyses.

6.4 Example output for the first 3 pairs run through the SSQuAWK3 workflow:

A	B	C	D	E	F	G	H	I	J	K
Sample	0xGenomeCov	<10xGenomeCov	nReads	avgLen	avgLenPassFilt	avgQual	avgQualPassFilt	avgCovPassQual	readsAlign	readsAlignPassFilt
SRR16828363.fastq.gz	107nt (0%)	31nt (0%)	632664	151	151	37.8	37.9	688X	138637 (21%)	136327 (21%)
SRR16828364.fastq.gz	76nt (0%)	31nt (0%)	458116	151	151	37.8	37.9	890X	179913 (39%)	176348 (38%)
SRR16828365.fastq.gz	76nt (0%)	77nt (0%)	351980	151	151	37.8	37.9	272X	54928 (15%)	53958 (15%)

Video Tutorial

7 Thanks for using SSQuAWK!



8 New to GalaxyTrakr? Check out this detailed, 12 minute video overview of the SSQuAWK (version 1) protocol before trying SSQuAWK3.

Video edit:

"SSQuawk allows users to check the sequence quality of SARS-CoV-2 wastewater samples in **CFSAN's custom Galaxy instance, called GalaxyTrakr**. This generates a single report file from raw Illumina MiSeq paired-end fastq file inputs."