

**TABLE OF CONTENTS – HYPERLINKS TO THE PROCEDURE**

- [Prepare a metadata file for the sequences to be uploaded to Terra \(5.1\)](#)
- [Log into Terra using Chrome and your Google account \(5.2\)](#)
- [Upload sequence files and metadata into Terra \(5.3\)](#)
- [Run the QC and genotyping workflow \(5.4\)](#)
- [Evaluate the QC metrics for the sequences \(5.5\)](#)
- [View the genotyping results for the sequences \(5.6\)](#)
- [Upload sequences to NCBI \(5.7\)](#)
- [Appendix PNID01-1: Data Import to Terra Directly from the Illumina BaseSpace](#)
- [Appendix PNID01-2: Data Download from NCBI SRA](#)
- [Appendix PNID01-3: Customization of a Data Table View for PulseNet QC Metrics](#)
- [Appendix PNID01-4a. PulseNet Critical Quality Metrics for Routine Sequence Submissions](#)
- [Appendix PNID01-4b. TheiaProk Read Pre-Screening Step to Exclude Poor Quality Sequences to Conserve Computational Resources](#)
- [Appendix PNID01-5. Customization of a Data Table View for PulseNet Genotyping Assays](#)
- [Appendix PNID01-6. Uploading Additional Metadata to Terra for NCBI Submissions and Customization of a Data Table View for Metadata](#)
- [Appendix PNID01-7: Overview of the TheiaProk Workflow for Bacterial Characterization](#)

**1. PURPOSE:** To describe the procedure for analyzing Illumina short read whole genome sequencing (WGS) data to be used for PulseNet International (PNI) surveillance utilizing the cloud-based Terra.Bio platform.

**2. SCOPE:** This procedure applies to all PulseNet personnel that utilize the Terra.Bio platform to analyze Illumina short read WGS data for surveillance activities within the PulseNet International network. This SOP covers sequence and metadata upload to Terra.Bio, sequence quality evaluation, assembly and genotyping workflows and sequence upload to NCBI. Phylogenetic analyses are covered by the SOP PNID02 (PulseNet International Standard Operating Procedure for Phylogenetic Analysis of WGS Data Using the Terra.Bio Platform).

**3. DEFINITIONS/TERMS:**

**3.1 ANI:** Average Nucleotide Intity

**3.2 BaseSpace:** Illumina cloud-based computing environment for next generation sequencing data analysis, management and storage, including data sharing.

**3.3 Bash Commands:** Bash (Bourne Again Shell) is a command-line interface (CLI) shell used extensively in Linux and macOS. Shell is a computer program that allows direct control of a computer’s operating system. Bash commands are used to control the computer or operating system without having to navigate menus, options, and windows within the graphical user interface.

- 3.4 BioProject:** A collection of biological data on NCBI related to a single initiative, originating from a single organization or from a consortium.
- 3.5 BioSample:** Submitter-supplied descriptive information (metadata) about the biological materials from which the NCBI stored data are derived.
- 3.6 Contig:** A contiguous consensus sequence derived from the assembly of many short, overlapping DNA fragments.
- 3.7 Coverage:** The average number of reads that includes a given nucleotide in the reconstructed sequence.
- 3.8 Critical Quality Metrics:** coverage (after trimming), average quality (Q score before trimming), assembly length, and secondary genus abundance (contamination detection by MIDAS). Sequences not meeting the minimum thresholds/acceptable ranges for these metrics defined in this document should be re-sequenced.
- 3.9 CSV:** Comma Separated Values
- 3.10 DeNovo assembly:** A sequence assembly generated from the short raw reads without the use of a reference genome.
- 3.11 FASTA:** A text-based format for representing either nucleotide sequences or peptide sequences, in which base pairs or amino acids are represented using single-letter codes. A sequence in FASTA format begins with a single-line description, followed by lines of sequence data. The description line is distinguished from the sequence data by a greater-than (" $>$ ") symbol in the first column.
- 3.12 FASTQ:** a text-based format for storing both a biological sequence and its corresponding quality scores.
- 3.13 GAMBIT:** Genomic Approximation Method for Bacterial Identification and Tracking. A bacterial species identification method utilizing k-mer-based algorithm to search against a large reference database of genomes.
- 3.14 Gzip:** A file format and a software application used for file compression and decompression to transfer data quickly over the Internet
- 3.15 LIMS:** Laboratory Information Management System
- 3.16 Mash Sketching:** Mash is a set of tools for creating and using MinHash sketches, a way of turning a genome into a small signature which can be easily compared with other signatures.
- 3.17 Metadata:** A set of data that describes and gives information about other data.
- 3.18 MIDAS:** Metagenomic Intra-species Diversity Analysis System. Integrated computational pipeline for quantifying bacterial species abundance and coverage from shotgun metagenomes based on blast alignment against a panel of universal single copy genes.
- 3.19 N50:** The N50 statistic is commonly used as a rough assessment of genomic assemblies. It represents the contig length (in base pairs) for which half of the genome sequence is assembled in contigs larger than or equal to N50 contig size.
- 3.20 NCBI:** National Center for Bio**te**chnology Information
- 3.21 PNI:** Pulse**N**et International
- 3.22 QA/QC:** Quality Assurance/Quality Control
- 3.23 Q score:** The quality score for each individual base position in a sequence, indicating the accuracy of the base call. Phred scores are used, where  $Q = -10\log(\text{Error Probability})$ . The



**PULSENET INTERNATIONAL STANDARD OPERATING PROCEDURE FOR ANALYZING ILLUMINA SHORT READ WGS DATA USING THE TERRA.BIO PLATFORM**

**Doc. No. PNID01**

**Ver. No. 01**

**Effective Date:**

**Page 4 of 61**

5.1.2 The following three columns are required at a minimum:

5.1.2.1 **Entity:collection\_name\_id**. This is the name of the Terra data collection (data table) into which you want to upload your sequences, e.g., entity:CDC\_ATCC\_Sequences\_id. The strain IDs (entry keys) are called “entities” in Terra.

**NOTE:** no spaces or dashes are allowed. “entity:” and “\_id” are required by Terra in the column name.

5.1.2.2 **Read1**. The name of the read 1 fastq.gz file for the entry.

5.1.2.3 **Read2**. The name of the read 2 fastq.gz file for the entry.

Entity	Read1	Read2
17802-C1-M947-23-007_S1_L001_R1_001.fastq.gz	17802-C1-M947-23-007_S1_L001_R1_001.fastq.gz	17802-C1-M947-23-007_S1_L001_R2_001.fastq.gz
17802-C2-M947-23-007_S2_L001_R1_001.fastq.gz	17802-C2-M947-23-007_S2_L001_R1_001.fastq.gz	17802-C2-M947-23-007_S2_L001_R2_001.fastq.gz
17802-C3-M947-23-007_S3_L001_R1_001.fastq.gz	17802-C3-M947-23-007_S3_L001_R1_001.fastq.gz	17802-C3-M947-23-007_S3_L001_R2_001.fastq.gz
33560-C1-M947-23-007_S4_L001_R1_001.fastq.gz	33560-C1-M947-23-007_S4_L001_R1_001.fastq.gz	33560-C1-M947-23-007_S4_L001_R2_001.fastq.gz
33560-C2-M947-23-007_S5_L001_R1_001.fastq.gz	33560-C2-M947-23-007_S5_L001_R1_001.fastq.gz	33560-C2-M947-23-007_S5_L001_R2_001.fastq.gz
33560-C3-M947-23-007_S6_L001_R1_001.fastq.gz	33560-C3-M947-23-007_S6_L001_R1_001.fastq.gz	33560-C3-M947-23-007_S6_L001_R2_001.fastq.gz
51812-C1-B-M947-23-007_S10_L001_R1_001.fastq.gz	51812-C1-B-M947-23-007_S10_L001_R1_001.fastq.gz	51812-C1-B-M947-23-007_S10_L001_R2_001.fastq.gz
51812-C2-B-M947-23-007_S11_L001_R1_001.fastq.gz	51812-C2-B-M947-23-007_S11_L001_R1_001.fastq.gz	51812-C2-B-M947-23-007_S11_L001_R2_001.fastq.gz
51812-C3-B-M947-23-007_S12_L001_R1_001.fastq.gz	51812-C3-B-M947-23-007_S12_L001_R1_001.fastq.gz	51812-C3-B-M947-23-007_S12_L001_R2_001.fastq.gz

**NOTE:** additional metadata can be added at this point or at later date. Refer to [appendix PNID01-6](#) for guidance on metadata.

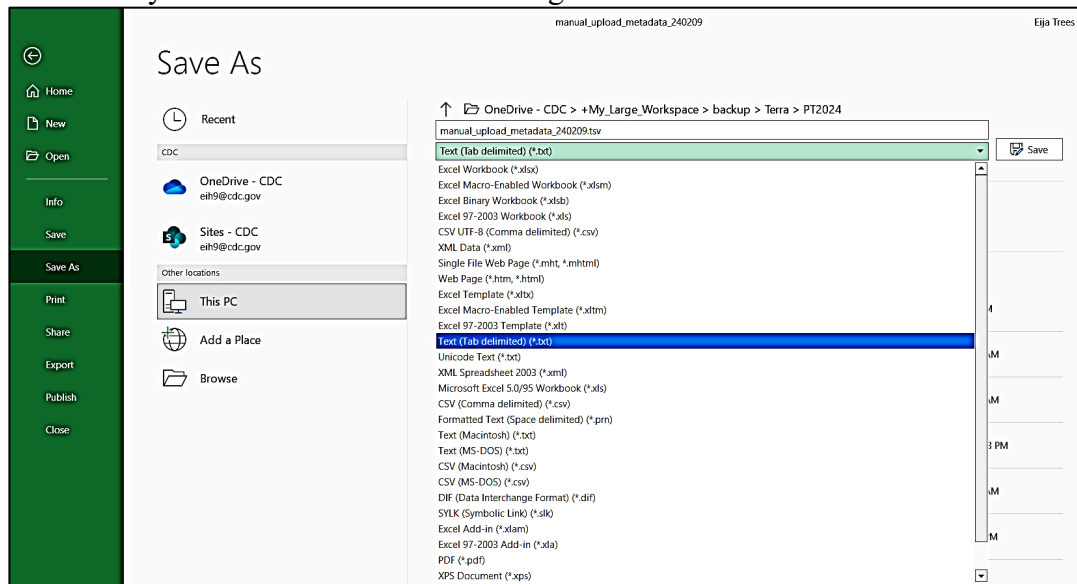
5.1.3 Type in the strain IDs in the “entity” column the way you want them to appear in Terra.

**NOTE:** the strain ID does not need to match any part of the fastq.gz file name.

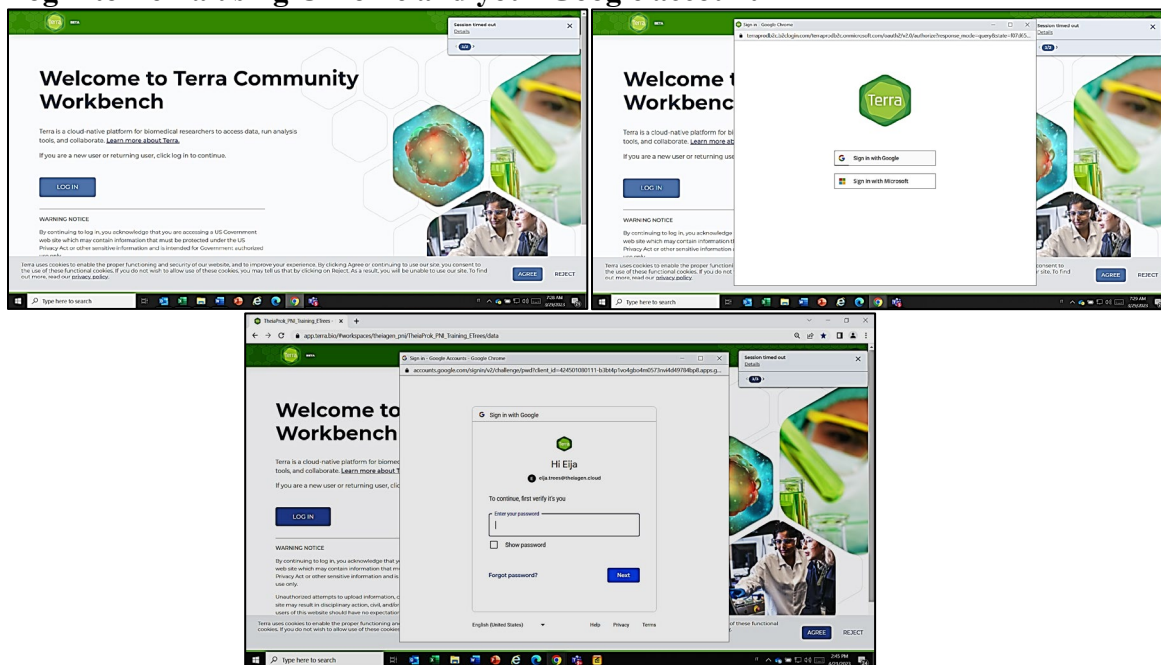
5.1.4 Copy and paste the fastq.gz file names for each strain in the columns “read1” and “read2”. Make sure the file names end “fastq.gz”.

Name	Date modified	Type	Size
New folder	2/8/2023 12:12 PM	File folder	
10708-C1-M3235-23-002_S16_L001_R1_001.fastq	1/25/2023 9:11 AM	GZ File	152,935 KB
10708-C1-M3235-23-002_S16_L001_R2_001.fastq	1/25/2023 9:11 AM	GZ File	175,618 KB
10708-C2-M3235-23-002_S17_L001_R1_001.fastq	1/25/2023 9:11 AM	GZ File	251,050 KB
10708-C2-M3235-23-002_S17_L001_R2_001.fastq	1/25/2023 9:11 AM	GZ File	273,935 KB
10708-C3-M3235-23-002_S18_L001_R1_001.fastq	1/25/2023 9:11 AM	GZ File	261,478 KB
10708-C3-M3235-23-002_S18_L001_R2_001.fastq	1/25/2023 9:12 AM	GZ File	293,655 KB
17802-C1-M947-23-007_S1_L001_R1_001.fastq		File	148,748 KB
17802-C1-M947-23-007_S1_L001_R2_001.fastq		File	128,117 KB
17802-C2-M947-23-007_S2_L001_R1_001.fastq		File	192,254 KB
17802-C2-M947-23-007_S2_L001_R2_001.fastq		File	197,471 KB
17802-C3-M947-23-007_S3_L001_R1_001.fastq		File	178,324 KB
17802-C3-M947-23-007_S3_L001_R2_001.fastq		File	188,470 KB
25922-C1-M3235-23-002_S19_L001_R1_001.fastq		File	290,392 KB
25922-C1-M3235-23-002_S19_L001_R2_001.fastq		File	312,502 KB
25922-C2-M3235-23-002_S20_L001_R1_001.fastq		File	208,370 KB
25922-C2-M3235-23-002_S20_L001_R2_001.fastq		File	224,490 KB
25922-C3-M3235-23-002_S21_L001_R1_001.fastq		File	243,397 KB
25922-C3-M3235-23-002_S21_L001_R2_001.fastq		File	277,029 KB
33560-C1-M947-23-007_S4_L001_R1_001.fastq	2/8/2023 4:41 PM	GZ File	806,679 KB
33560-C1-M947-23-007_S4_L001_R2_001.fastq	2/8/2023 4:41 PM	GZ File	99,987 KB

5.1.5 Save the file in the **tsv format**: choose “Save As” and “Text (Tab delimited) (\*.txt)”. Make sure your file name has **“.tsv”** ending.



## 5.2 Log into Terra using Chrome and your Google account



## 5.3 Upload sequence files and metadata into Terra

5.3.1 Under the “Terra Workspaces”, select the “Data” tab, click “Import data” and select “Open data uploader” from the drop-down menu.

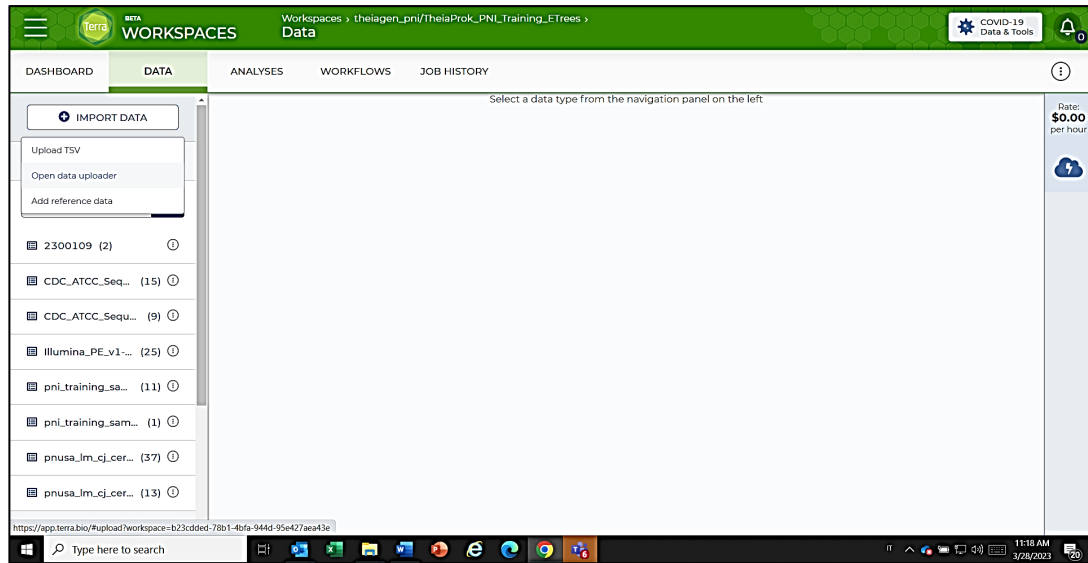
**PULSENET INTERNATIONAL STANDARD OPERATING PROCEDURE FOR ANALYZING ILLUMINA SHORT READ WGS DATA USING THE TERRA.BIO PLATFORM**

**Doc. No. PNID01**

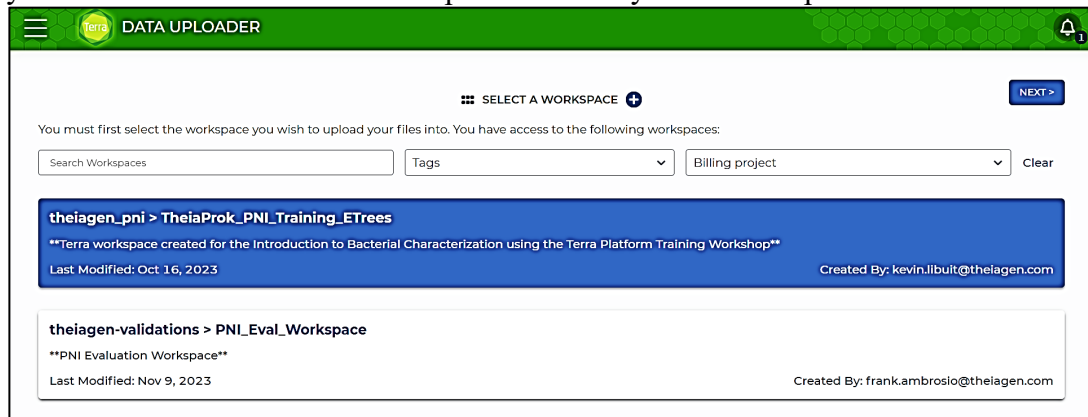
**Ver. No. 01**

**Effective Date:**

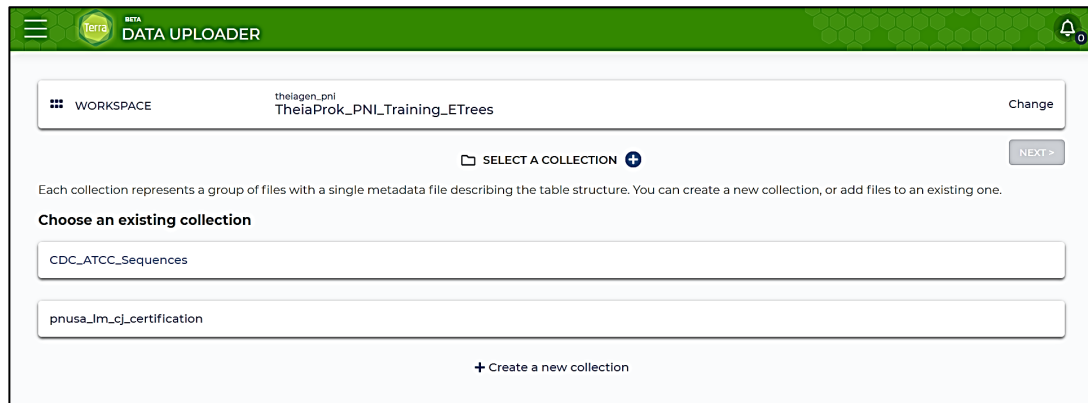
**Page 6 of 61**



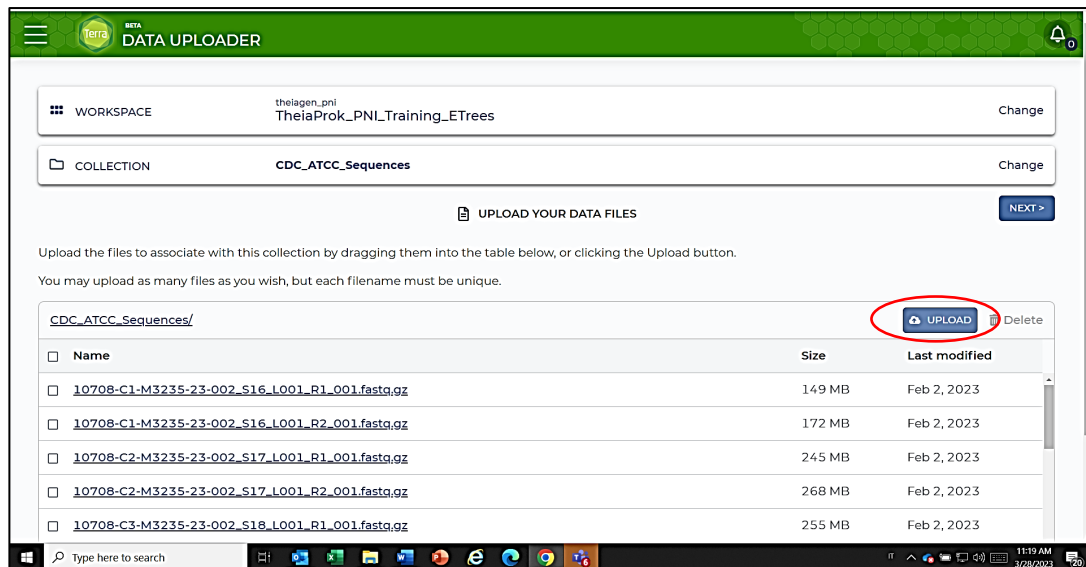
5.3.2 “Data uploader” screen will open. If your account has access to multiple workspaces, you first need to choose the workspace to which you wish to upload the data.



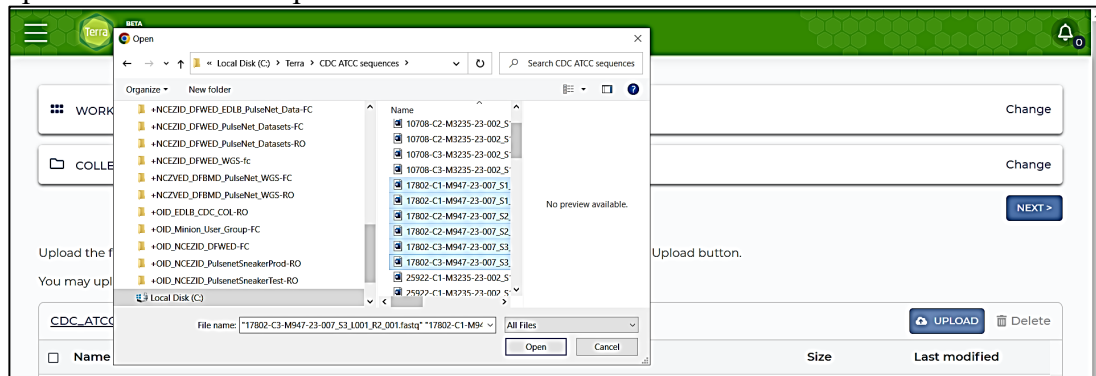
5.3.3 Either choose an existing collection by clicking on the collection name on the list or create a new collection.



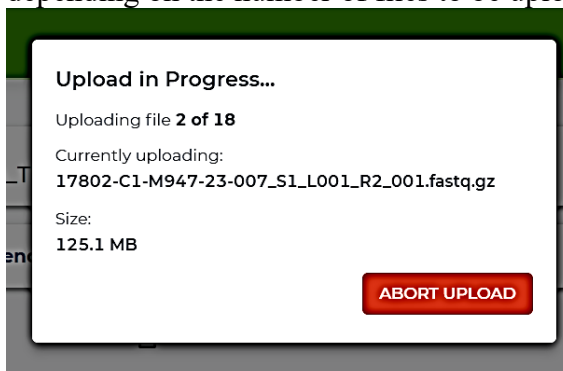
5.3.4 Under “Upload your data files”, click on “Upload”.



5.3.5 Navigate to the location where the FASTQ files are saved, select the files to be uploaded and click “Open”.



5.3.6 “Upload in progress” pop-up window will appear. The upload may take a few minutes depending on the number of files to be uploaded and your Internet bandwidth.



5.3.7 After the pop-up window disappears, click on “Next”.

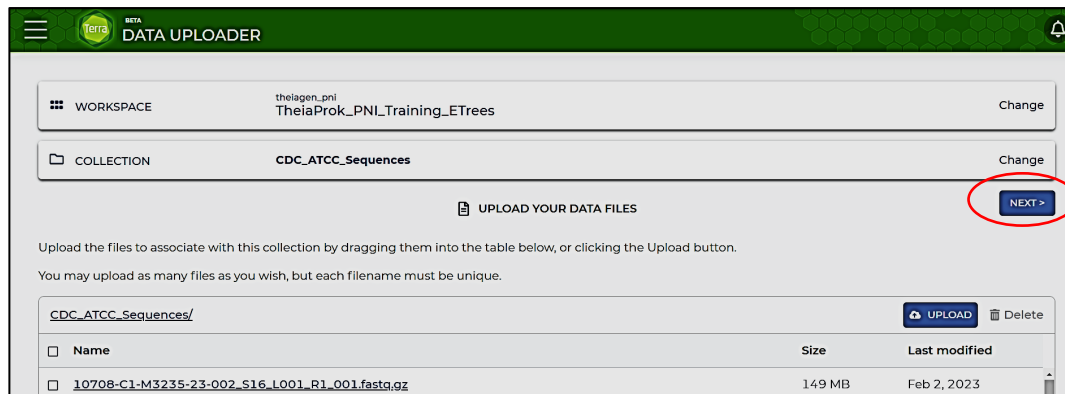
# PULSENET INTERNATIONAL STANDARD OPERATING PROCEDURE FOR ANALYZING ILLUMINA SHORT READ WGS DATA USING THE TERRA.BIO PLATFORM

Doc. No. PNID01

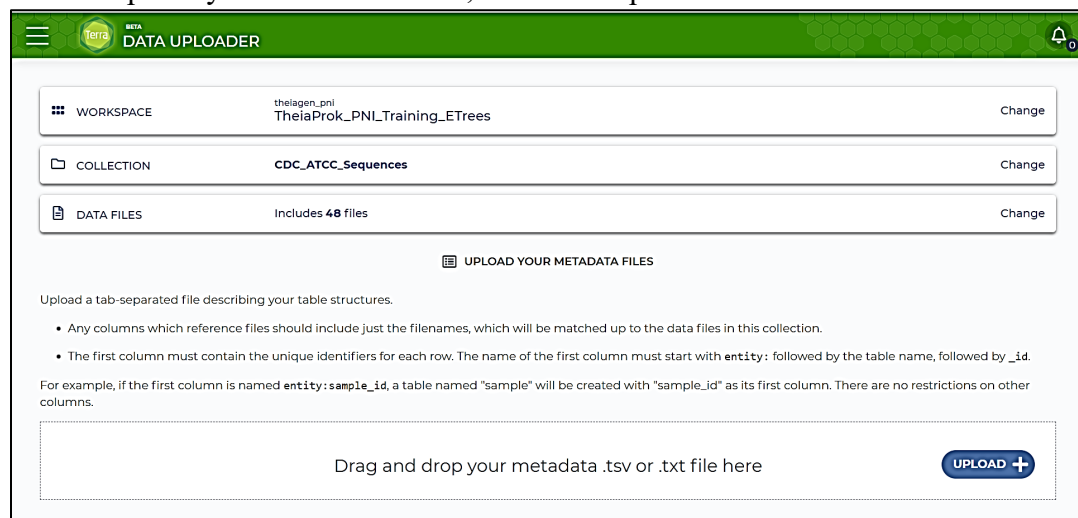
Ver. No. 01

Effective Date:

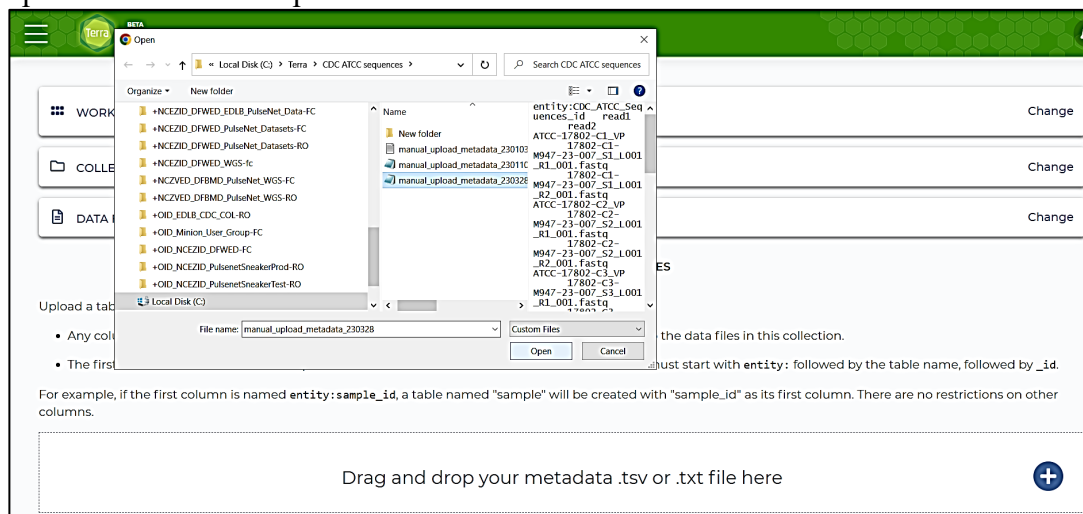
Page 8 of 61



5.3.8 Under “Upload your metadata files”, click on “Upload”.



5.3.9 Navigate to the location where the metadata tsv file is saved, select the file to be uploaded and click “Open”.



5.3.10 On the following screen, check that the fastq.gz files are linked to the correct entry keys and click “Update table”.

**PULSENET INTERNATIONAL STANDARD OPERATING PROCEDURE FOR ANALYZING ILLUMINA SHORT READ WGS DATA USING THE TERRA.BIO PLATFORM**

**Doc. No. PNID01**

**Ver. No. 01**

**Effective Date:**

**Page 9 of 61**

Updating Table: **CDC\_ATCC-Sequences** RENAME TABLE CANCEL **UPDATE TABLE**

▲ This workspace already includes a table with this name. If any new rows have the same CDC\_ATCC-Sequences\_id as an existing row, the data in that row will be updated with the new values.  
If this table looks right to you, click the button on the right to update the table in your workspace.

entity: CDC_ATCC_Sequence...	▲ read1 (updated)	▲ read2 (updated)
ATCC-17802-C1_VP	17802-C1-M947-23-007_S1_L001_R1_001.fastq.gz	17802-C1-M947-23-007_S1_L001_R2_001.fastq.gz
ATCC-17802-C2_VP	17802-C2-M947-23-007_S2_L001_R1_001.fastq.gz	17802-C2-M947-23-007_S2_L001_R2_001.fastq.gz
ATCC-17802-C3_VP	17802-C3-M947-23-007_S3_L001_R1_001.fastq.gz	17802-C3-M947-23-007_S3_L001_R2_001.fastq.gz
ATCC-33560-C1_CJ	33560-C1-M947-23-007_S4_L001_R1_001.fastq.gz	33560-C1-M947-23-007_S4_L001_R2_001.fastq.gz
ATCC-33560-C2_CJ	33560-C2-M947-23-007_S5_L001_R1_001.fastq.gz	33560-C2-M947-23-007_S5_L001_R2_001.fastq.gz
ATCC-33560-C3_CJ	33560-C3-M947-23-007_S6_L001_R1_001.fastq.gz	33560-C3-M947-23-007_S6_L001_R2_001.fastq.gz
ATCC-51812-C1_SE	51812-C1-B-M947-23-007_S10_L001_R1_001.fastq...	51812-C1-B-M947-23-007_S10_L001_R2_001.fastq.gz
ATCC-51812-C2_SE	51812-C2-B-M947-23-007_S11_L001_R1_001.fastq...	51812-C2-B-M947-23-007_S11_L001_R2_001.fastq.gz
ATCC-51812-C3_SE	51812-C3-B-M947-23-007_S12_L001_R1_001.fastq...	51812-C3-B-M947-23-007_S12_L001_R2_001.fastq.gz

5.3.11 A “Done” message will appear on the “Data uploader” screen. You can view the updated data table by clicking on the link that appears on the screen.

**NOTE:** *The columns visible in the data table can be customized. For the ease of navigation, it is recommended to create separate views for QC metrics, genotyping results and metadata. Refer to appendices [PNID01-3](#) (QC metrics), [PNID01-5](#) (Genotyping) and [PNID01-6](#) (Metadata) for guidance on how to customize the data table columns for PulseNet surveillance.*

Terra BETA DATA UPLOADER

WORKSPACE **theiagen\_pni**  
TheiaProk\_PNI\_Training\_ETrees Change

COLLECTION **CDC\_ATCC-Sequences** Change

DATA FILES **Includes 48 files** Change

METADATA TABLES **Updated table CDC\_ATCC-Sequences, added or modified 9 rows** Change

**DONE!**

- View the CDC\_ATCC-Sequences table in the workspace ←
- Create a new table in the CDC\_ATCC-Sequences collection
- Start over with another workspace or collection

**PULSENET INTERNATIONAL STANDARD OPERATING PROCEDURE FOR ANALYZING ILLUMINA SHORT READ WGS DATA USING THE TERRA.BIO PLATFORM**

Doc. No. PNID01

Ver. No. 01

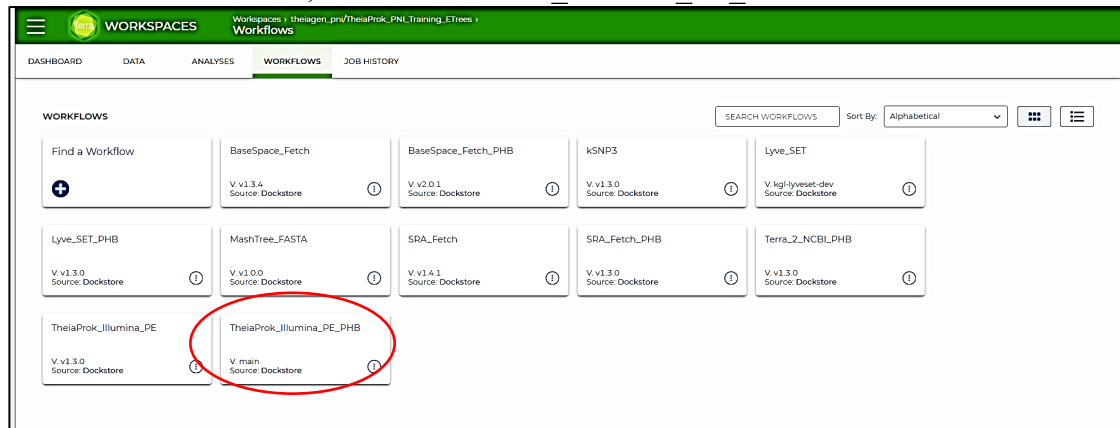
Effective Date:

Page 10 of 61

CDC_ATCC_Sequences_id	number_co...	read1	read2
<input type="checkbox"/> ATCC-10708-C1_SE	51	10708-C1-M3235-23-002_S16_L001_R1_001.fastq.gz	10708-...
<input type="checkbox"/> ATCC-10708-C2_SE	49	10708-C2-M3235-23-002_S17_L001_R1_001.fastq.gz	10708-...
<input type="checkbox"/> ATCC-10708-C3_SE	50	10708-C3-M3235-23-002_S18_L001_R1_001.fastq.gz	10708-...
<input type="checkbox"/> ATCC-17802-C1_VP		17802-C1-M947-23-007_S1_L001_R1_001.fastq.gz	17802-...
<input type="checkbox"/> ATCC-17802-C2_VP		17802-C2-M947-23-007_S2_L001_R1_001.fastq.gz	17802-...
<input type="checkbox"/> ATCC-17802-C3_VP		17802-C3-M947-23-007_S3_L001_R1_001.fastq.gz	17802-...
<input type="checkbox"/> ATCC-17802_VP	48	ATCC-17802-M947-22-040_S4_L001_R1_001.fastq.gz	ATCC-1...
<input type="checkbox"/> ATCC-25922-C1_EC	60	25922-C1-M3235-23-002_S19_L001_R1_001.fastq.gz	25922-...
<input type="checkbox"/> ATCC-25922-C2_EC	68	25922-C2-M3235-23-002_S20_L001_R1_001.fastq.gz	25922-...

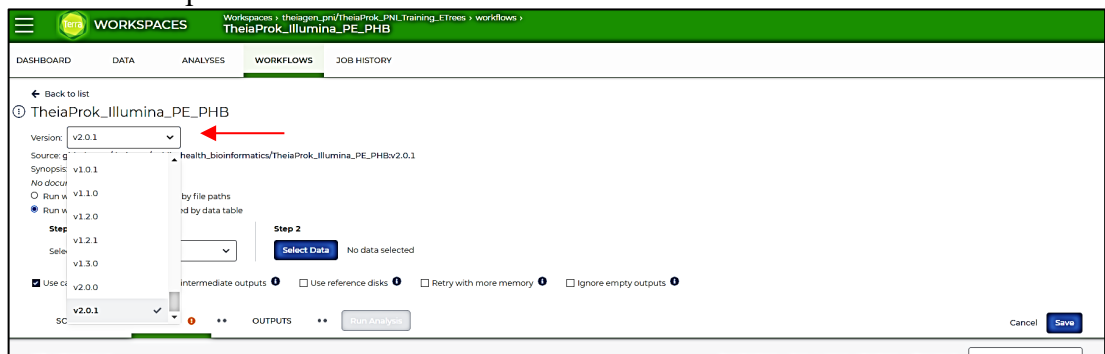
5.4 Run the QC and genotyping workflow. The TheiaProk workflow performs the QC of the raw sequence reads, de novo assembles the raw reads using Skesa and then performs the QC of the assembly and identification of the species. A variety of genotyping assays appropriate to the species are also available.

5.4.1 In the “Workflows” tab, select “TheiaProk\_Illumina\_PE\_PHB”.



5.4.2 In the “TheiaProk\_Illumina\_PE\_PHB” screen:

5.4.2.1 Select the latest version of the TheiaProk\_Illumina\_PE\_PHB workflow from the “Version” drop-down menu.



5.4.2.2 Under “Step 1”, select the “Root type”, i.e., the data table in which the samples are located, e.g., “CDC\_ATCC\_Sequences”.

**NOTE:** For most data tables, there are two options: the **main** data table containing the individual sample entries and the **set** data table containing sample sets used for phylogenetic analyses, NCBI uploads, etc (e.g., CDC\_ATCC\_Sequences\_set). Make sure to select the **main** data table.

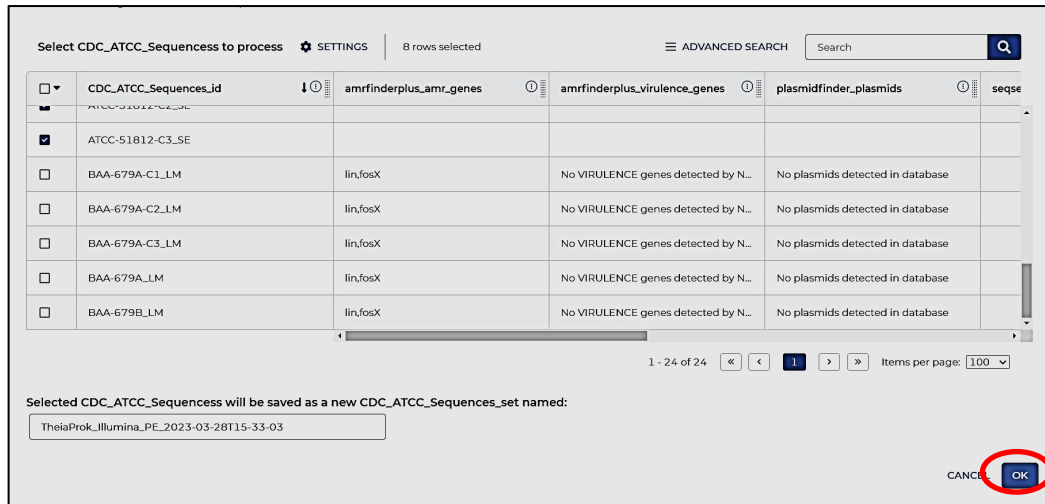
Variable	Type	Input value
read1	File	Required
read2	File	Required
samplename	String	Required

5.4.2.3 Under “Step 2”, click “Select data” (screenshot above). This will take you to the data table specified in Step 1.

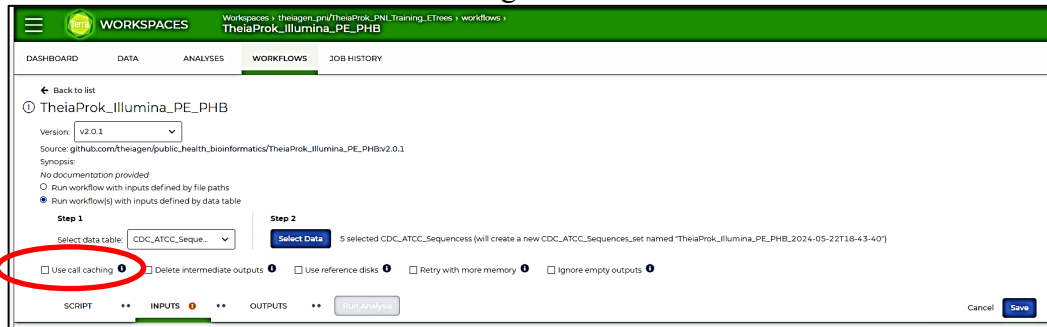
5.4.2.4 Select the strains to be analyzed and scroll all the way to the bottom to click “OK”.

**NOTE:** if the data table contains more than 100 entries and you check the “Select all” box by the data table name, only 100 entries will be selected.

CDC_ATCC_Sequences_id	amrfinderplus_amr_genes	amrfinderplus_virulence_genes	plasmidfinder_plasmids	seqse
<input type="checkbox"/> ATCC-10708-C1_SE	mdsA,mdsB	sinH,iroB,iroC,sodC1	IncFIB(S),IncFII(S)	7:c:1...
<input type="checkbox"/> ATCC-10708-C2_SE	mdsB,mdsA	sinH,iroB,iroC,sodC1	IncFIB(S),IncFII(S)	7:c:1...
<input type="checkbox"/> ATCC-10708-C3_SE	mdsB,mdsA	sinH,iroB,iroC,sodC1	IncFIB(S),IncFII(S)	7:c:1...
<input checked="" type="checkbox"/> ATCC-17802-C1_VP				
<input checked="" type="checkbox"/> ATCC-17802-C2_VP				
<input checked="" type="checkbox"/> ATCC-17802-C3_VP				
<input type="checkbox"/> ATCC-17802_VP				



5.4.2.5 De-select the box for “Use call catching”.



5.4.2.6 In the “Inputs” tab, specify the following input values in the “Attribute” column (scroll down the list):

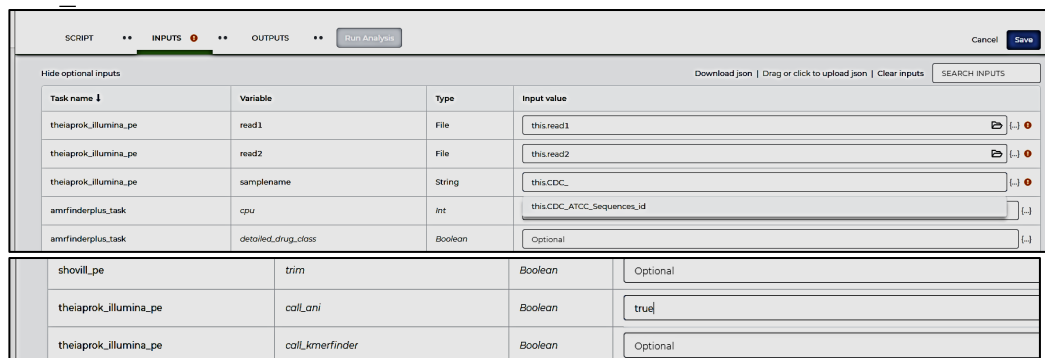
**NOTE:** When you fill in the Attribute column, clicking inside the cell will bring up a drop-down menu of attributes that you can select to avoid typos (screenshot below).

5.4.2.6.1 Read1: “This.read1”.

5.4.2.6.2 Read2: “This.read2”.

5.4.2.6.3 Samplename: “**This.data table name\_id**”, e.g., This.CDC\_ATCC\_Sequences\_id.

5.4.2.6.4 Call\_ani: “True”.



5.4.2.7 In the “Outputs” tab, click “Use defaults”.

**PULSENET INTERNATIONAL STANDARD OPERATING PROCEDURE FOR ANALYZING ILLUMINA SHORT READ WGS DATA USING THE TERRA.BIO PLATFORM**

**Doc. No. PNID01**

**Ver. No. 01**

**Effective Date:**

**Page 13 of 61**

Task name	Variable	Type	Input value   Use defaults
theiaprok_illumina_pe	abricate_abau_plasmid_tsv	File	this.abricate_abau_plasmid_tsv
theiaprok_illumina_pe	abricate_abau_plasmid_type_genes	String	this.abricate_abau_plasmid_type_genes
theiaprok_illumina_pe	abricate_database	String	this.abricate_database
theiaprok_illumina_pe	abricate_docker	String	this.abricate_docker
theiaprok_illumina_pe	abricate_version	String	this.abricate_version
theiaprok_illumina_pe	agrivate_agr_canonical	String	this.agrivate_agr_canonical
theiaprok_illumina_pe	agrivate_agr_group	String	this.agrivate_agr_group
theiaprok_illumina_pe	agrivate_agr_match_score	String	this.agrivate_agr_match_score

5.4.2.8 Click “Save” (screenshot above).

**NOTE:** the “Save” button is only visible if you have changed the inputs from the previous submission.

5.4.2.9 Click “Run analysis”. A “Confirm launch” pop-up window appears that allows you to type in an optional description. Click “Launch”.

**Confirm launch**

Output files will be saved as workspace data in: us US (multi-region)

Running workflows will generate cloud charges

How much does my workflow cost? ☑  
Set up budget alert ☑

Describe your submission (optional):  
Colony picks for VP, C3 and ST

This will launch 8 analyses.

CANCEL LAUNCH

5.4.3 A “Workflow statuses” screen will appear where your submitted jobs should be initially listed as “Queued”.

Data Entity	Last Changed	Status	Run Cost	Messages	Workflow ID
ATCC-10708-C1_SE (CDC_ATCC-Sequences)	May 23, 2024, 8:08 AM	Queued	N/A		
ATCC-10708-C2_SE (CDC_ATCC-Sequences)	May 23, 2024, 8:08 AM	Queued	N/A		
ATCC-10708-C3_SE (CDC_ATCC-Sequences)	May 23, 2024, 8:08 AM	Queued	N/A		

5.4.4 Go to the “Job history” tab to check the status of your job submission. A successfully finished job is indicated by a green check mark.

# PULSENET INTERNATIONAL STANDARD OPERATING PROCEDURE FOR ANALYZING ILLUMINA SHORT READ WGS DATA USING THE TERRA.BIO PLATFORM

Doc. No. PNID01

Ver. No. 01

Effective Date:

Page 14 of 61

Submission (click for details)	Data entity	No of Workflows	Status	Submitted ↑	Submission ID	Comment	Actions
TheiaProk_illumina_PE_PHB Submitted by ejaj.trees@theiagen.cloud	TheiaProk_illumina_PE_PHB...	3	✓ Done	May 23, 2024 8:09 AM	a25e9acc-935ee-4097-ac72- e075d6a08632	Repeat of the choleraesuis cert strain colony...	ⓘ
BaseSpace_Fetch_PHB Submitted by ejaj.trees@theiagen.cloud	BaseSpace_Fetch_PHB_202...	3	✓ Done	May 7, 2024 1:45 PM	37f7e097-71ac-4abc-950b- 86487371c0d9	£ additional sequences from the CAOC Bas...	ⓘ
BaseSpace_Fetch_PHB Submitted by ejaj.trees@theiagen.cloud	BaseSpace_Fetch_PHB_202...	10	▲ Done	May 6, 2024 3:42 PM	bf858564-1bb0-4f69-8718- e95438940bdc	10 samples from the BaseSpace nextseq nu...	ⓘ
BaseSpace_Fetch_PHB_Ix5pDqZMIA Submitted by ejaj.trees@theiagen.cloud	BaseSpace_Fetch_PHB_202...	10	✓ Done	May 6, 2024 1:45 PM	952544fb-6298-43a3-8932- b1237506e699	10 additional samples from BaseSpace next...	ⓘ
TheiaProk_illumina_PE Submitted by ejaj.trees@theiagen.cloud	TheiaProk_illumina_PE_202...	10	✓ Done	May 6, 2024 7:52 AM	87807612-3957-4ac6-a363- cdc0e49595ec	10 samples from v3 validation run VL403-2...	ⓘ

5.5 Evaluate the QC metrics for the sequences: QC metrics can be viewed either directly on the data table (5.5.1-5.5.3) or they can be exported to Excel for the selected entries (5.5.4).

5.5.1 Under the “Terra Workspaces”, select the “Data” tab, then select the data table of interest, e.g., “CDC\_ATCC\_Sequences”.

TABLES	EDIT	OPEN WITH...	EXPORT	SETTINGS	0 rows selected	ADVANCED SEARCH	SEARCH
<input type="checkbox"/> CDC_ATCC_Sequences_id <input type="checkbox"/> ATCC-10708-C1_SE <input type="checkbox"/> ATCC-10708-C2_SE <input type="checkbox"/> ATCC-10708-C3_SE <input type="checkbox"/> ATCC-17802-C1_VP <input type="checkbox"/> ATCC-17802-C2_VP <input type="checkbox"/> ATCC-17802-C3_VP <input type="checkbox"/> ATCC-17802_VP <input type="checkbox"/> ATCC-25922-C1_EC <input type="checkbox"/> ATCC-25922-C2_EC <input type="checkbox"/> ATCC-25922-C3_EC							

5.5.2 Select “Settings”.

TABLES	EDIT	OPEN WITH...	EXPORT	SETTINGS	0 rows selected	ADVANCED SEARCH	SEARCH
<input type="checkbox"/> CDC_ATCC_Sequences_id <input type="checkbox"/> ATCC-10708-C1_SE <input type="checkbox"/> ATCC-10708-C2_SE <input type="checkbox"/> ATCC-10708-C3_SE <input type="checkbox"/> ATCC-17802-C1_VP				Change the order and visibility of columns in the table			

5.5.3 On the “Select columns” screen under “Your saved column selections”, click on the circle with 3 dots next to the “qc\_metrics” and from the drop-down menu select “Load” and then click “Done”. This will load the appropriate PulseNet QC metrics into the data table. Refer to the Appendix [PNID01-3](#) for the QC metrics that are supposed to appear on the table and for the instructions on how to add or delete any of the columns (QC metrics) in the QC metrics table.



# PULSENET INTERNATIONAL STANDARD OPERATING PROCEDURE FOR ANALYZING ILLUMINA SHORT READ WGS DATA USING THE TERRA.BIO PLATFORM

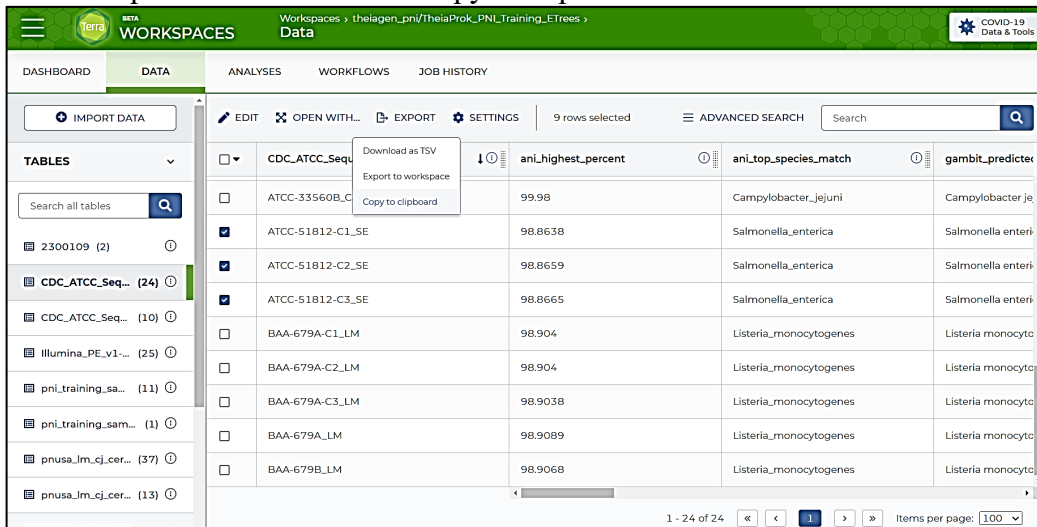
Doc. No. PNID01

Ver. No. 01

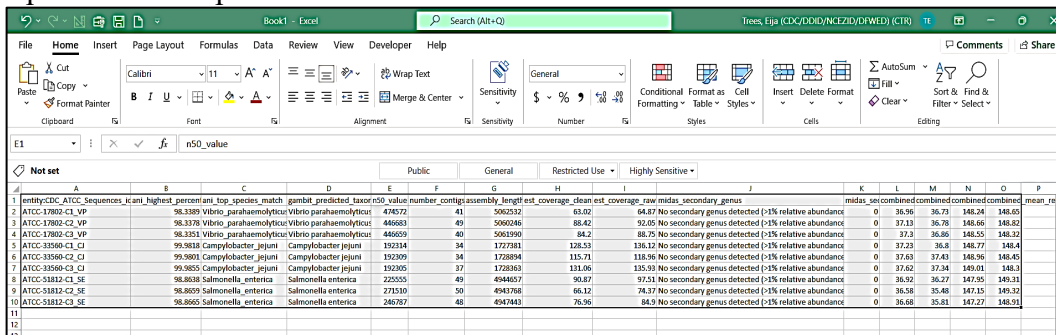
Effective Date:

Page 16 of 61

## 5.5.4.2 Click “Export” and then select “Copy to clipboard”.

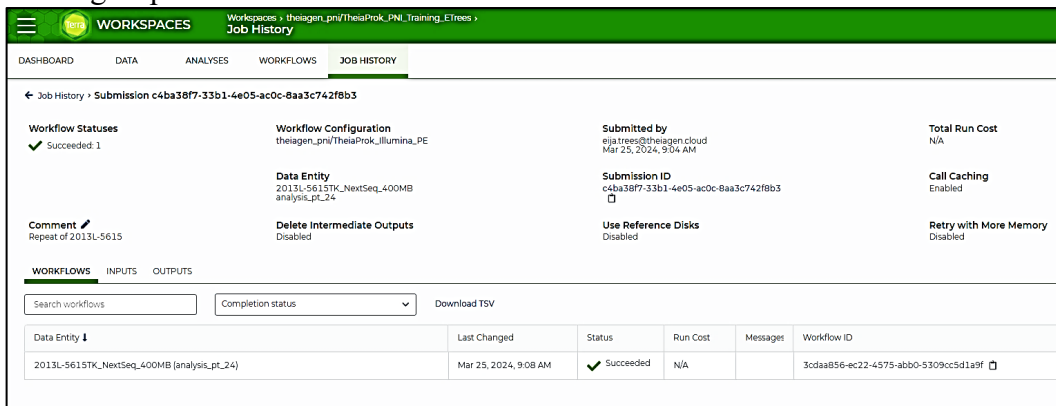


## 5.5.4.3 Open Excel and paste the data to a worksheet.



5.5.5 Refer to Appendix [PNID01-4a](#) for PulseNet critical quality metrics for acceptable Illumina sequences.

5.5.6 TheiaProk employs a read pre-screening step in which it will stop analysis for samples that fall below certain quality thresholds. In this case, the “Job History” tab indicates that the sample was successful, but there are no results in the “Data” tab. The “raw\_read\_screen” column in qc metrics should indicate the reason for the sequence analysis failure. Refer to Appendix [PNID01-4b](#) for thresholds applied in this read pre-screening step.



**PULSENET INTERNATIONAL STANDARD OPERATING PROCEDURE FOR ANALYZING ILLUMINA SHORT READ WGS DATA USING THE TERRA.BIO PLATFORM**

**Doc. No. PNID01**

**Ver. No. 01**

**Effective Date:**

**Page 17 of 61**

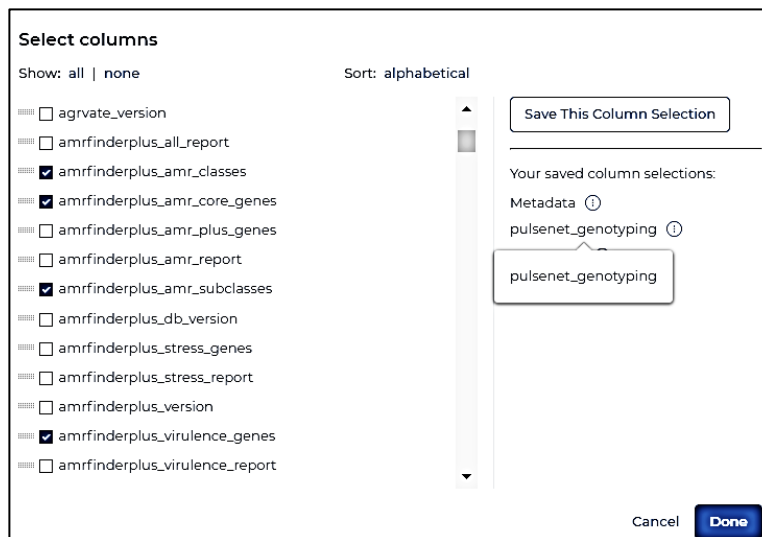
analysis_pt_24_id	midas_secondary_genus	midas_secondary_genus_abunda...	n50_value	number_sortigs	raw_read_screen
2011V-1043_FLEX_300_Vbrno	No secondary genus detected (>1% r...	0	124349	77	PASS
2012V-1116_FLEX_300_Vbrno	No secondary genus detected (>1% r...	0	458045	39	PASS
2013L-536L_FLEX_300_LM	No secondary genus detected (>1% r...	0.0008	526928	12	PASS
2013L-5410_FLEX_300_LM	No secondary genus detected (>1% r...	0.0008	526025	15	PASS
2013L-5547_FLEX_300_LM	No secondary genus detected (>1% r...	0	435363	20	PASS
2013L-5615TK_NextSeq_400MB					FAIL; the estimated coverage is less than the minimum of 10x
2015AM-1304	No secondary genus detected (>1% r...	0	443204	15	PASS
2015AM-1305	No secondary genus detected (>1% r...	0	728098	16	PASS

5.6 **View the genotyping results for the sequences:** Genotyping results can be viewed either directly on the data table (5.6.1-5.6.3) or they can be exported to Excel for the selected entries (5.6.4).

5.6.1 Under the “Terra Workspaces”, select the “Data” tab, then select the data table of interest, e.g., “CDC\_ATCC\_Sequences”.

5.6.2 In the “Data” tab, select “Settings”.

5.6.3 On the “Select columns” screen under “Your saved column selections”, click on the circle with 3 dots next to the “pulsenet\_genotyping” and from the drop-down menu select “Load” and then click “Done”. This will load the genotyping assays appropriate for the PulseNet surveillance into the data table. Refer to the Appendix [PNID01-5](#) for the genotyping assays that are supposed to appear on the table and for the instructions on how to add or delete any of the columns (genotyping assays) in the genotyping results table.



5.6.4 Export the results to Excel for the selected entries; follow the procedure in the step 5.5.4.

**5.7 Upload sequences to NCBI**

**NOTE:** Contact Theiagen Genomics ([support@theiagen.com](mailto:support@theiagen.com)) for guidance before starting and to configure your workspace for the NCBI uploads. The configuration process is described in: [https://theiagen.notion.site/Terra\\_2\\_NCBI-61abcdc066646b3b258f70b561e9f62](https://theiagen.notion.site/Terra_2_NCBI-61abcdc066646b3b258f70b561e9f62).

**PULSENET INTERNATIONAL STANDARD OPERATING PROCEDURE FOR ANALYZING ILLUMINA SHORT READ WGS DATA USING THE TERRA.BIO PLATFORM**

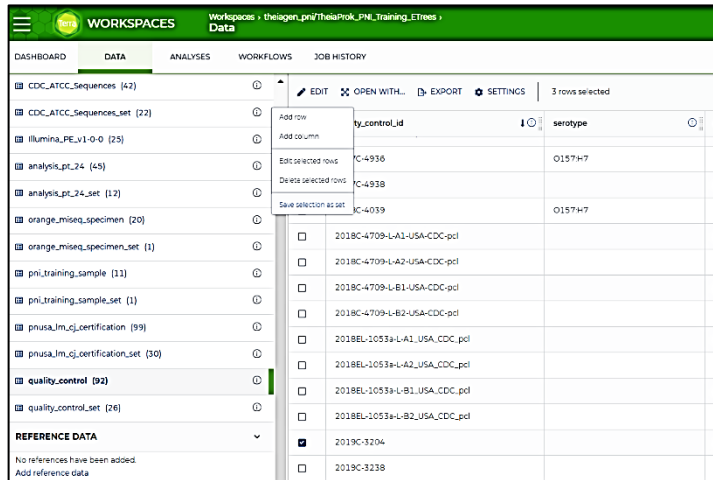
**Doc. No. PNID01**

**Ver. No. 01**

**Effective Date:**

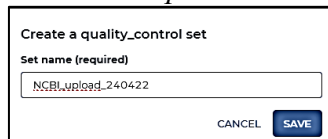
**Page 18 of 61**

- 5.7.1 Upload the metadata required for the NCBI submission: refer to the [appendix PNID01-6](#) for the correct metadata formatting and upload.
- 5.7.2 Create a **set** of the samples to be uploaded to NCBI:
  - 5.7.2.1 Under the “Terra Workspaces”, select the “Data” tab, then select the data table of interest, e.g., “quality\_control”.
  - 5.7.2.2 In the “Data” tab, select the sequences to be included in the NCBI upload.
  - 5.7.2.3 From the “Edit” drop-down menu, select “Save selection as set”.

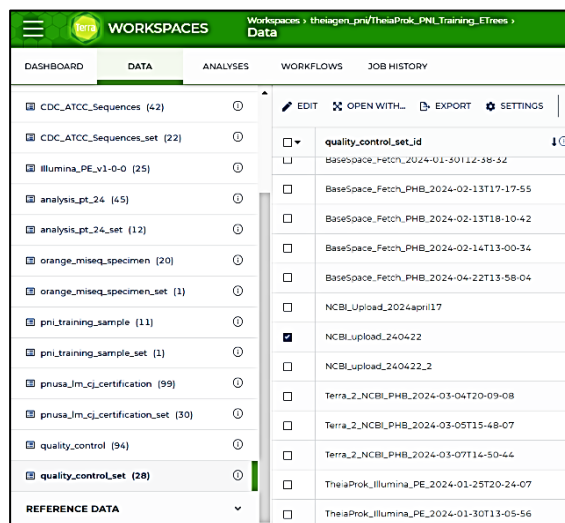


- 5.7.2.4 In the appearing pop-up window, name the set, e.g., “NCBI\_upload\_240422” and click “Save”.

**NOTE:** *no spaces or dashes are allowed.*

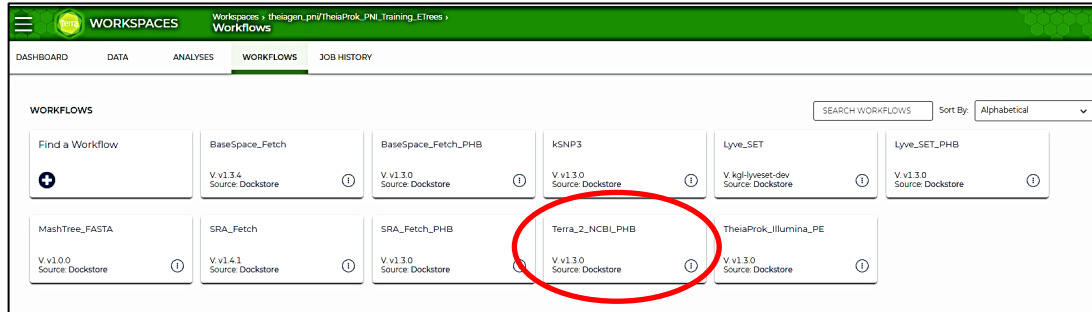


- 5.7.2.5 The newly created set should now appear in the “set” data table, e.g., “quality\_control\_set”.

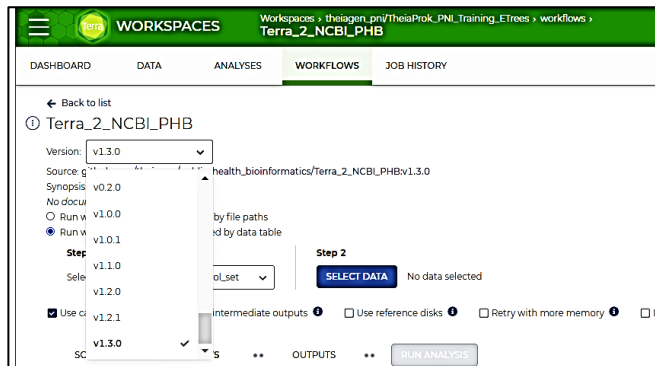


5.7.3 Set up parameters for the NCBI upload workflow:

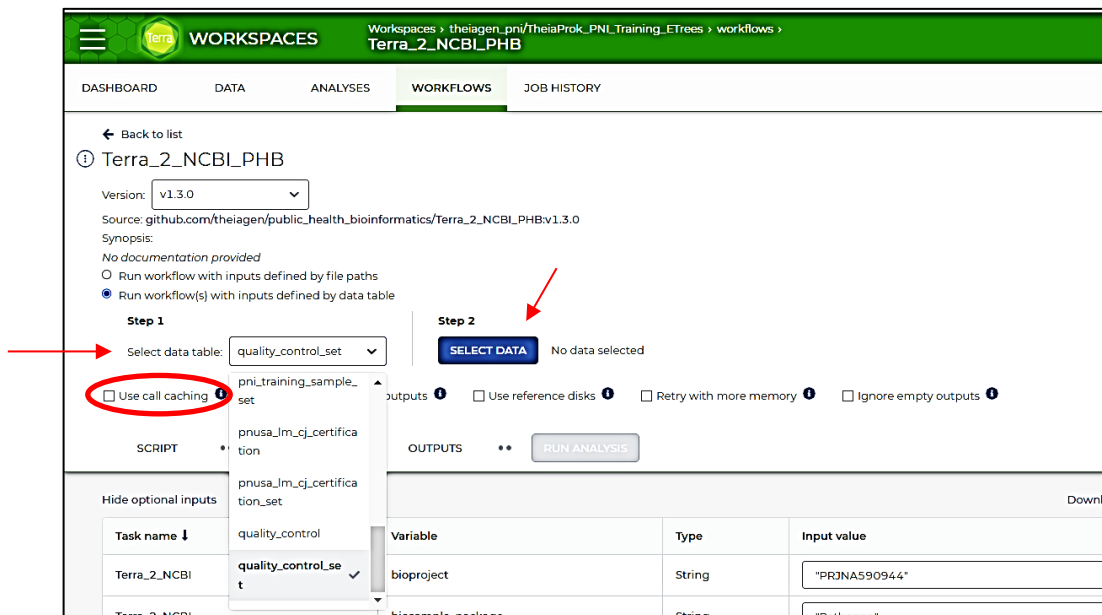
5.7.3.1 In the “Workflows” tab, select the “Terra\_2\_NCBI\_PHB” workflow. This will bring up the “Terra\_2\_NCBI\_PHB” screen.



5.7.3.2 From the “Version” drop-down menu, select the latest version of the “Terra\_2\_NCBI\_PHB”.

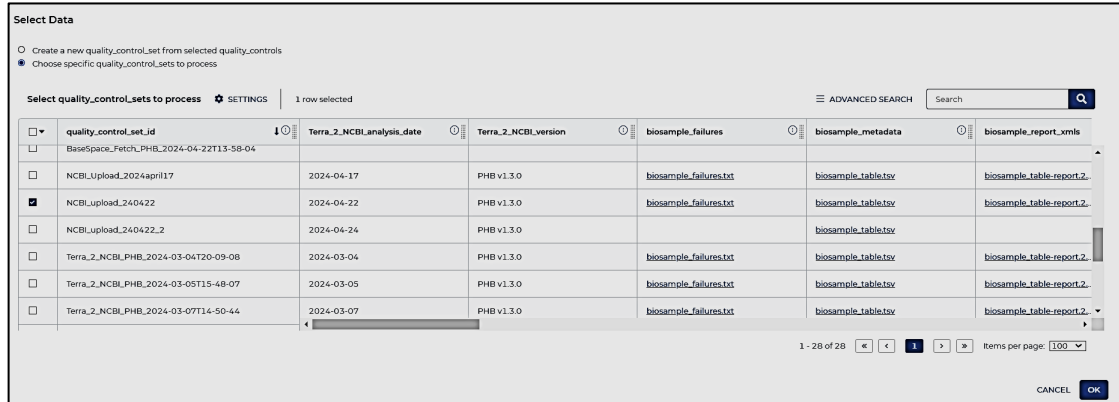


5.7.3.3 Under Step 1, from the “Select root entity type” drop-down menu, select the set data table where the sample set created in step 5.7.2. is located, e.g., “quality\_control\_set”. Also de-select the box for “Use call catching”.



5.7.3.4 Under Step 2, click “Select data” (screenshot above). This will take you to the set data table selected in the previous step.

5.7.3.5 Select the desired sample set, e.g., “NCBI\_upload\_240422” and click “OK”.



5.7.3.6 In the “Inputs” tab, the following “Input Values” need to be filled out for the “Variables” listed below:

5.7.3.6.1 bioproject: enter the number for the NCBI BioProject to which you wish to upload in quotation marks, e.g., “PRJNA590944”.

5.7.3.6.2 biosample\_package in quotation marks: “Pathogen”. This is the metadata template/package you use for the metadata upload for sequences belonging to PulseNet surveillance.

5.7.3.6.3 ncbi\_config\_js: enter the name of the NCBI configuration file created for your workspace, e.g., workspace.ncbi\_config\_etrees.

5.7.3.6.4 project\_name in quotation marks: “theiagen\_pni”.

5.7.3.6.5 sample\_names: enter your data table name in the format:

**this.data table names.data table name\_id,**  
e.g., this.quality\_controls.quality\_control\_id.

**NOTE: the double name format is REQUIRED.**

5.7.3.6.6 sra\_transfer\_gcp\_bucket in quotation marks: “gs://theiagen\_sra\_transfer”. This is the temporary public Google storage location for your sequences that NCBI can access.

5.7.3.6.7 table\_name: enter your data table name in quotation marks, e.g., “quality\_control”.

5.7.3.6.8 workspace\_name: enter your workspace name in quotation marks e.g., “TheiaProk\_PNI\_Training\_ETrees”.

5.7.3.6.9 submit\_to\_production: true.

5.7.3.6.10 OPTIONAL: if you want to link an SRA submission to an **existing BioSample** (“biosample\_accession” field already populated with the SAMN number in the “Data” tab):

5.7.3.6.10.1 skip\_biosample: true.

**NOTE: for using this feature, the required metadata fields need to be populated, i.e., populating only the “biosample\_accession” field is not enough and will not pass the Terra pre-upload checks.**

**PULSENET INTERNATIONAL STANDARD OPERATING PROCEDURE FOR ANALYZING ILLUMINA SHORT READ WGS DATA USING THE TERRA.BIO PLATFORM**

**Doc. No. PNID01**

**Ver. No. 01**

**Effective Date:**

**Page 21 of 61**

Task name ↓	Variable	Type	Input value
Terra_2_NCBI	bioproject	String	"PRJNA590944"
Terra_2_NCBI	biosample_package	String	"Pathogen"
Terra_2_NCBI	ncbi_config_js	File	workspace ncbi_config_etrees
Terra_2_NCBI	project_name	String	"theiagen_pni"
Terra_2_NCBI	sample_names	Array[String]	this.quality_controls.quality_control_id
Terra_2_NCBI	sra_transfer_gcp_bucket	String	"gs://theiagen_sra_transfer"
Terra_2_NCBI	table_name	String	"quality_control"
Terra_2_NCBI	workspace_name	String	"TheiaProk_PNI_Training_ETrees"
ncbi_sftp_upload	additional_files	Array[File]	Optional
ncbi_sftp_upload	wait_for	String	Optional
prune_table	read1_column_name	String	Optional
prune_table	read2_column_name	String	Optional
Terra_2_NCBI	input_table	File	Optional
Terra_2_NCBI	skip_biosample	Boolean	Optional
Terra_2_NCBI	submit_to_production	Boolean	true

5.7.3.7 In the “**Outputs**” tab, click “Use defaults” for the “Attributes” and then click “Save”.  
**NOTE:** the “Save” button is only visible if you have changed the inputs from the previous submission.

Task name ↓	Variable	Type	Input value   Use defaults
Terra_2_NCBI	biosample_failures	File	this.biosample_failures
Terra_2_NCBI	biosample_metadata	File	this.biosample_metadata
Terra_2_NCBI	biosample_report_xmils	Array[File]	this.biosample_report_xmils
Terra_2_NCBI	biosample_status	String	this.biosample_status
Terra_2_NCBI	biosample_submission_xml	File	this.biosample_submission_xml
Terra_2_NCBI	excluded_samples	File	this.excluded_samples
Terra_2_NCBI	generated_accessions	File	this.generated_accessions
Terra_2_NCBI	sra_metadata	File	this.sra_metadata

5.7.3.8 Click “Run analysis” (screenshot above). A “Confirm launch” pop-up window appears that allows you to type in an optional description. Click “Launch”.

**Confirm launch**

Output files will be saved as workspace data in:  
us-us-central1 (Iowa) ⓘ

Running workflows will generate cloud charges. ⓘ  
How much does my workflow cost? ⓘ  
Set up budget alert ⓘ

Describe your submission (optional):

NCBI upload of 3 sequences to the validation  
 bioproject

This will launch 1 analysis.

CANCEL
LAUNCH

**PULSENET INTERNATIONAL STANDARD OPERATING PROCEDURE FOR ANALYZING ILLUMINA SHORT READ WGS DATA USING THE TERRA.BIO PLATFORM**

**Doc. No. PNID01**

**Ver. No. 01**

**Effective Date:**

**Page 22 of 61**

5.7.3.9 A “Workflow statuses” screen will appear where your submitted jobs should be initially listed as “Queued” or “Launching”.

The screenshot shows the 'Job History' page for a submission with ID 387be0c7-88b6-4c00-8028-9020dc5f08b8. The workflow is in a 'Submitted' state with a clock icon. The configuration is 'theigen\_pnl/terra\_2\_NCBL\_PHB'. The data entity is 'NCBLupload240422 quality\_control\_set'. The submission was made by 'eja.trees@theigen.cloud' on Apr 22, 2024, at 7:47 AM. The total run cost is N/A. The comment indicates 'NCBI upload of 3 sequences to the validation...'. The 'Delete Intermediate Outputs' and 'Use Reference Disks' options are disabled. A table below shows the workflow details:

Data Entity	Last Changed	Status	Run Cost	Messages	Workflow ID
NCBLupload240422 (quality_control_set)	Apr 22, 2024, 7:47 AM	Queued	N/A		

5.7.3.10 Go to the “Job history” tab to check the status of your job submission. A successfully finished job is indicated by a green check mark.

The screenshot shows a list of job submissions in the 'Job History' tab. The table below summarizes the data:

Submission (click for details)	Data entity	Nd. of Workflows	Status	Submitted	Submission ID	Comment	Actions
Terra_2_NCBL_PHB Submitted by eja.trees@theigen.cloud	NCBLupload240422 (quality...	1	Done	Apr 22, 2024 7:47 AM	387be0c7-88b6-4c00-8028-9020dc5f08b8	NCBI upload of 3 sequences to the...	
Terra_2_NCBL_PHB Submitted by eja.trees@theigen.cloud	NCBLupload2024a0117 (...	1	Done	Apr 17, 2024 7:53 AM	8aac979-2965-4549-a06a-436c9f9d327	NCBI submission of 3 new sample...	
TheiaProk_illumina_PE Submitted by eja.trees@theigen.cloud	TheiaProk_illumina_PE_202...	6	Done	Apr 15, 2024 12:22 PM	4bb20e4-f8b-4942-80c-516850a818a4	6 strains from the AMD incubator ...	
TheiaProk_illumina_PE Submitted by eja.trees@theigen.cloud	TheiaProk_illumina_PE_202...	3	Done	Apr 15, 2024 1:52 PM	5d443bc-51ba-4d94-9e3b-4673640f7627	Colony pick 2 for 2023 PI strains	
TheiaProk_illumina_PE Submitted by eja.trees@theigen.cloud	TheiaProk_illumina_PE_202...	3	Done	Apr 15, 2024 12:02 PM	0817f6c1-dc1f-4a43-a377-4408eaf1068	3 differentiated from 2023	

5.7.3.11 Tracking the NCBI accession numbers:

5.7.3.11.1 BioSample accession numbers (SAMN numbers)

5.7.3.11.1.1 In the “Data” tab, go to the data table in question, and you should see the “biosample\_accession” field populated with the SAMN number that is assigned to each unique BioSample.

The screenshot shows the 'Data' tab with a table of data entities. The 'quality\_control\_id' column is selected, and the 'biosample\_accession' column is highlighted with a red circle. The table contains the following data:

quality_control_id	biosample_accession	collected_by
2017C-4936	SAMN41039458	CDC
2017C-4938		
2018C-4039	SAMN41039457	CDC
2018C-4709-L-A1-USA-CDC-pcl		
2018C-4709-L-A2-USA-CDC-pcl		
2018C-4709-L-B1-USA-CDC-pcl		
2018C-4709-L-B2-USA-CDC-pcl		
2018EL-1053a-L-A1_USA_CDC_pcl		
2018EL-1053a-L-A2_USA_CDC_pcl		
2018EL-1053a-L-B1_USA_CDC_pcl		
2018EL-1053a-L-B2_USA_CDC_pcl		
2019C-3204	SAMN41039456	CDC
2019C-3238		

**PULSENET INTERNATIONAL STANDARD OPERATING PROCEDURE FOR ANALYZING ILLUMINA SHORT READ WGS DATA USING THE TERRA.BIO PLATFORM**

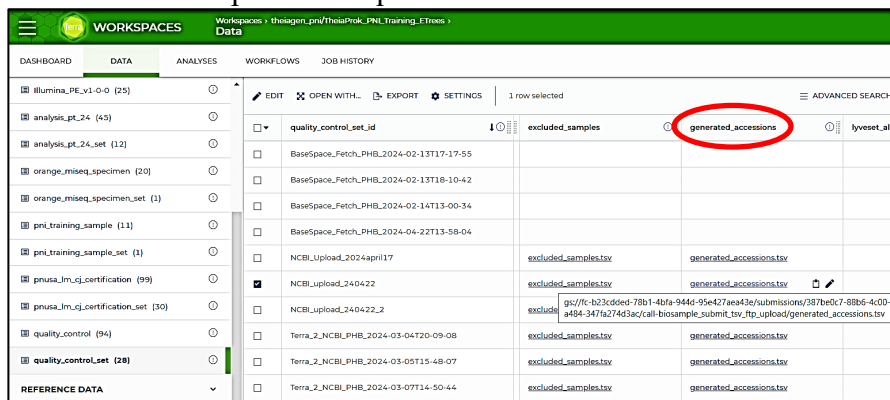
**Doc. No. PNID01**

**Ver. No. 01**

**Effective Date:**

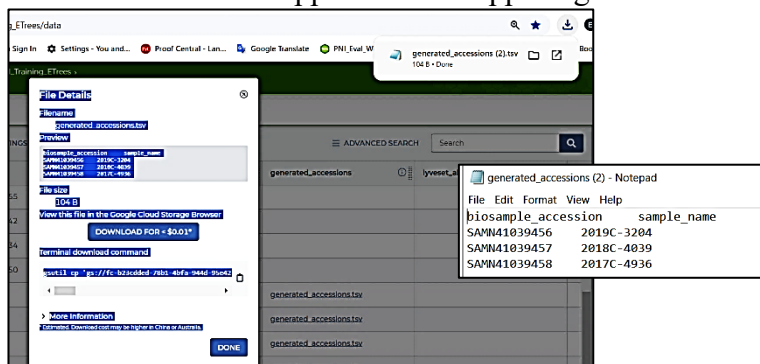
**Page 23 of 61**

5.7.3.11.1.2 In the “Data” tab, go to the data table **set** in question, find the data set you created for the NCBI upload and then click on the “generated\_accessions” link for the tsv file that lists the BioSample numbers for the uploaded sequence set. To download the file:

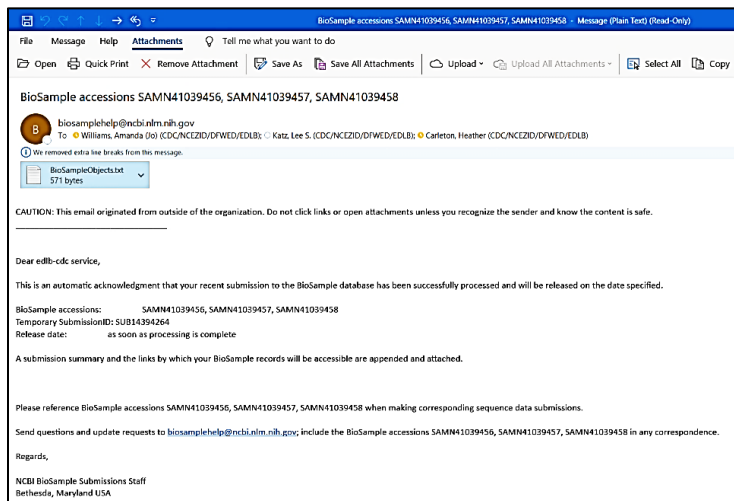


5.7.3.11.1.2.1 Click “Download >\$0.01\*” and then click “Done”.

5.7.3.11.1.2.2 The downloaded file appears on the upper right corner.



5.7.3.11.1.3 The email address associated with the NCBI account used for the upload should receive a confirmation email listing the SAMN numbers in the body of the email and in the “BiosampleObjects” text file attached to the email.



**PULSENET INTERNATIONAL STANDARD OPERATING PROCEDURE FOR ANALYZING ILLUMINA SHORT READ WGS DATA USING THE TERRA.BIO PLATFORM**

**Doc. No. PNID01**

**Ver. No. 01**

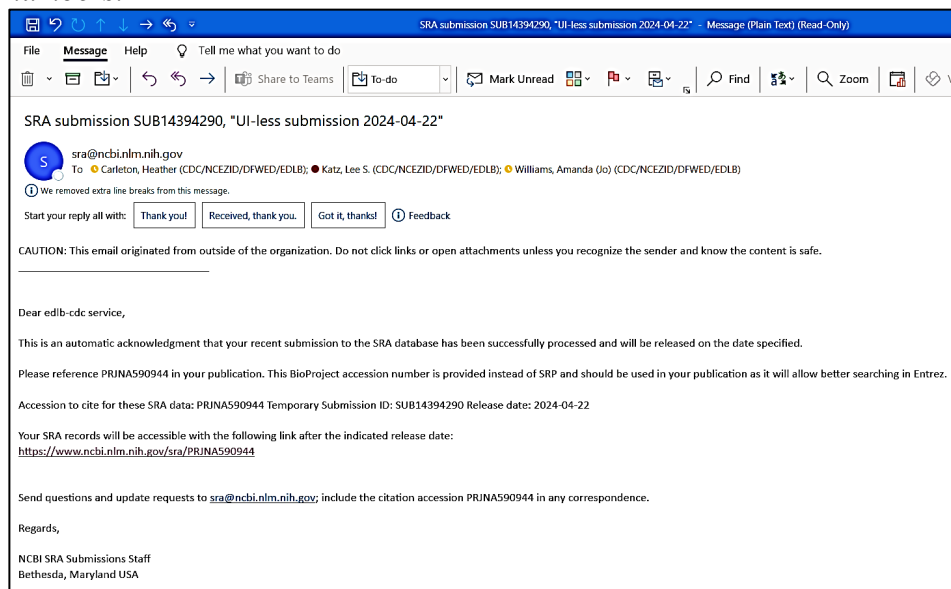
**Effective Date:**

**Page 24 of 61**

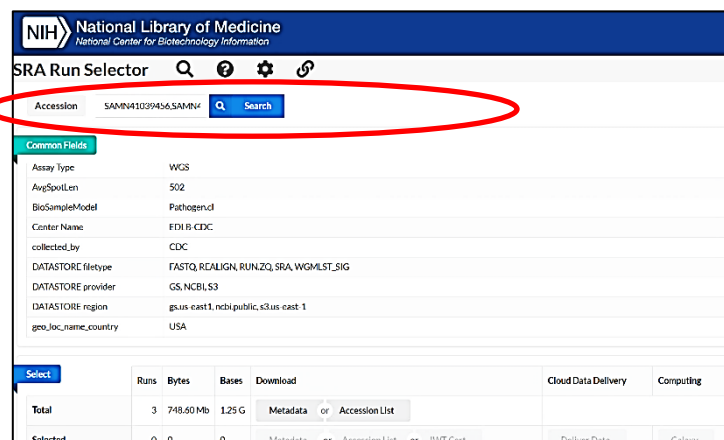
Accession	SPUID	Organism	Tax ID	Strain	BioProject
SAMN41039458	2017C-4936	Escherichia coli	562	2017C-4936	PRJNA590944
SAMN41039457	2018C-4039	Escherichia coli	562	2018C-4039	PRJNA590944
SAMN41039456	2019C-3204	Shigella sonnei 624	2019C-3204	2019C-3204	PRJNA590944

**5.7.3.11.2 SRA accession numbers (SRR numbers)**

**NOTE:** *The email address associated with the NCBI account used for the upload should receive a second confirmation email stating that the submission to the SRA database has been successfully processed and will be released on the date specified in the email. The email does not contain the SRA accession numbers.*



**5.7.3.11.2.1** Copy and paste the SAMN numbers from the tsv file from step 5.7.3.11.1.2. or the txt file from step 5.7.3.11.1.3. to the “NCBI Run Selector” tool at: <https://0-www-ncbi-nlm-nih-gov.brum.beds.ac.uk/Traces/study/>. Separate the numbers with commas. Click on “Search”.



**PULSENET INTERNATIONAL STANDARD OPERATING PROCEDURE FOR ANALYZING ILLUMINA SHORT READ WGS DATA USING THE TERRA.BIO PLATFORM**

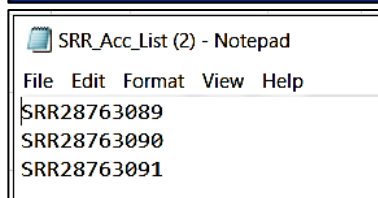
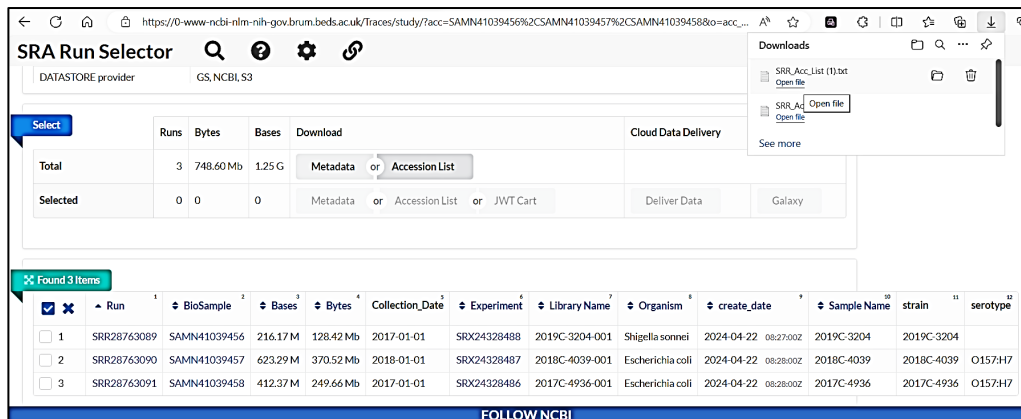
**Doc. No. PNID01**

**Ver. No. 01**

**Effective Date:**

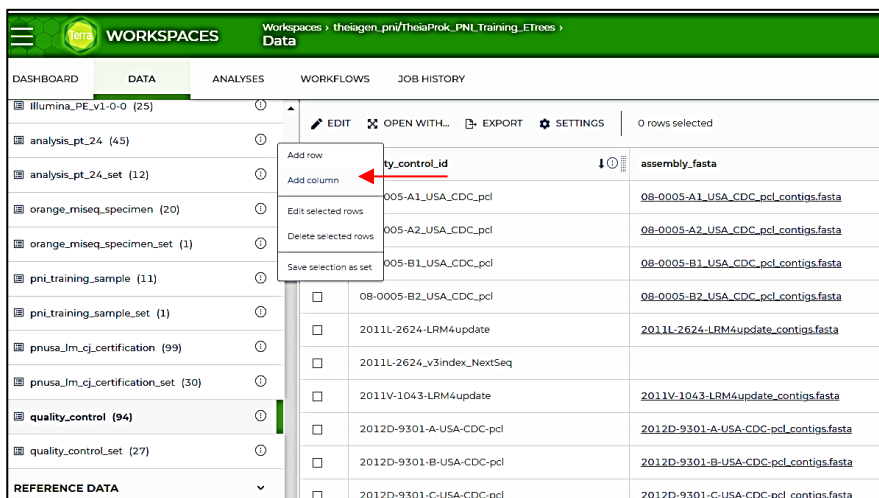
**Page 25 of 61**

- 5.7.3.11.2.2 The search will return a table that contains all the information NCBI has about your sequences, including the SRR accession numbers.
- 5.7.3.11.2.3 Click on “Accession List” to download the table containing the accession numbers. The downloaded file will appear on the upper right corner.



- 5.7.3.11.2.4 You can track the SRA accession numbers either on a separate Excel spreadsheet, LIMS system or you can create an “sra\_accession” field in your Terra data table and copy and paste the accession numbers there for each sample entry. To create a sra\_accession column and use it for tracking:

5.7.3.11.2.4.1 From the “Edit” drop-down menu, select “Add column”.



- 5.7.3.11.2.4.2 In the “Add a new column” pop-up window, name the new column “sra\_accession” and click “Save”.

**Add a new column**

**Column name**

**Default value** (optional, will be entered for all rows)

Type:

String    Reference    Number    Boolean

Value is a list

Cancel   **Save**

5.7.3.11.2.4.3 The new column should appear in the data table. To enter a SRR accession number for a specific sequence, click on “Edit value” in the sra\_accession column for that sample.

quality_control_id	sra_accession	assembly_fasta
2015K-0887_v3index_NextSeq		
2015K-1104_v3index_NextSeq		
2015K-1440-LRM4update		
2017C-3818-LRM4update		2017C-3818-LRM4update_contigs
2017C-3830-LRM4update		
2017C-4936		2017C-4936_contigs.fasta
2017C-4938		2017C-4938_contigs.fasta
2018C-4039		2018C-4039_contigs.fasta
2018C-4709-L-A1-USA-CDC-pcl		2018C-4709-L-A1-USA-CDC-pcl
2018C-4709-L-A2-USA-CDC-pcl		2018C-4709-L-A2-USA-CDC-pcl

5.7.3.11.2.4.4 Paste the SRA accession number to the field in the “Edit value” pop-up window and click “Save changes”.

**Edit value**

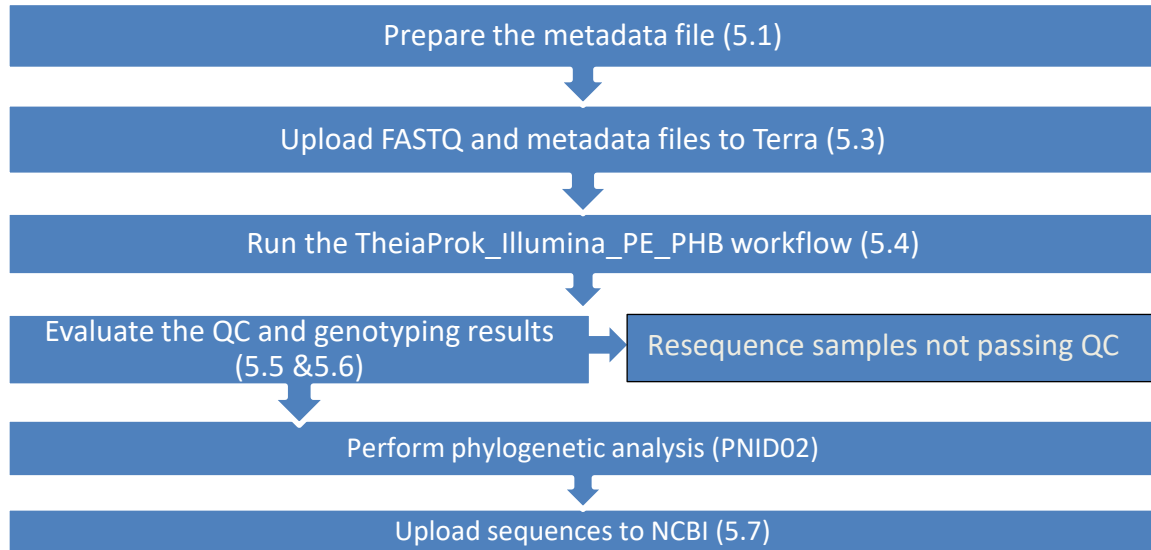
Type:

String    Reference    Number    Boolean

Value is a list

DELETE   CANCEL   **SAVE CHANGES**

## 6. FLOW CHART:



## 7. RELATED DOCUMENTS:

- 7.1 **PNID02:** PulseNet International Standard Operating Procedure for Phylogenetic Analysis of WGS Data Using the Terra.Bio Platform.

## 8. REFERENCES:

- 8.1 Libuit K.G., Doughty E.L., Otieno J.R., Ambrosio F., Kapsak C.J., Smith E.A., Wright S.M., Scribner M.R., Petit III R.A., Mendes C.I., Huergo M., Legacki G., Loreth C., Park D.J., Sevinsky J.R. (2023) Accelerating bioinformatics implementation in public health. *Microbial Genomics* 9:001051.

## 9. CONTACTS:

- 9.1 CDC USA PulseNet NGS Laboratory: [pulsenetngslab@cdc.gov](mailto:pulsenetngslab@cdc.gov)  
9.2 PulseNet International Quality Assurance Coordinator Eija Trees: [ehyytia-trees@cdc.gov](mailto:ehyytia-trees@cdc.gov)  
9.3 Theiagen:  
9.3.1 Generic email for support : [support@theiagen.com](mailto:support@theiagen.com)  
9.3.2 Michelle Scribner: [michelle.scribner@theiagen.com](mailto:michelle.scribner@theiagen.com)  
9.3.3 Frank Ambrosio: [frank.ambrosio@theiagen.com](mailto:frank.ambrosio@theiagen.com)

## 10. AMENDMENTS: NA

**PULSENET INTERNATIONAL STANDARD OPERATING PROCEDURE FOR ANALYZING ILLUMINA SHORT READ WGS DATA USING THE TERRA.BIO PLATFORM**

**Doc. No. PNID01**

**Ver. No. 01**

**Effective Date:**

**Page 28 of 61**

**11. APPROVAL SIGNATURES:**

Approved By: \_\_\_\_\_ Date: \_\_\_\_\_  
PulseNet QA/QC Personnel

Approved By: \_\_\_\_\_ Date: \_\_\_\_\_  
PulseNet WGS Technical Lead

Approved By: \_\_\_\_\_ Date: \_\_\_\_\_  
PulseNet International Coordinator

Approved By: \_\_\_\_\_ Date: \_\_\_\_\_  
PulseNet Response and Outbreak Management Team Lead

Approved By: \_\_\_\_\_ Date: \_\_\_\_\_  
Enteric Diseases Laboratory Branch Chief

**Appendix PNID01-1: Data Import to Terra Directly from the Illumina BaseSpace**

**NOTE:** In order to configure your Terra workspace to connect to your Illumina BaseSpace account, follow the instructions found in the Theiagen resources site at [https://theiagen.notion.site/BaseSpace\\_Fetch-34978656aa2d46ba82f2059434bd9369](https://theiagen.notion.site/BaseSpace_Fetch-34978656aa2d46ba82f2059434bd9369). For further assistance, contact Theiagen (see the Contacts section (9) for contact information).

1. Log into your BaseSpace account and find the run to be imported to Terra.

STATUS	RUN NAME	AVX/NQ30	%PF	INSTRUMENT	CREATED
Complete	WV-M07896-231215	...	81.72%	48.18% M07896	2023-12-15 15:18
Complete	WV-M07896-231213	...	90.55%	82.75% M07896	2023-12-13 12:18
Complete	VL403-23-003	...	88.63%	79.83% VL00403	2023-12-12 13:35
Complete	WV-M07896-231128	...	46.99%	10.94% M07896	2023-11-28 16:40
Complete	<b>M3235-23-042</b>	...	92.90%	91.02% M03235	2023-11-07 16:16
Complete	M3235-23-041	...	88.00%	92.64% M03235	2023-11-03 12:01
Complete	IMR-M01432-241023	...	84.67%	69.95% M01432	2023-10-24 01:26
Complete	IMR-M01432-171023	...	84.35%	68.96% M01432	2023-10-17 01:05
Complete	OH-VH01632-230919	...	89.77%	76.79% VH01632	2023-09-19 16:25
Complete	MiSeq Vc 015	...	91.99%	94.84% M08444	2023-09-01 15:45

2. Download the SampleSheet for the run:
  - a. Switch to the “Files” tab and scroll all the way to the bottom.

**M3235-23-042**

SUMMARY BIOSAMPLES SAMPLES CHARTS METRICS INDEXING QC SAMPLE SHEET **FILES**

Instrument: M03235 (92.90% NQ30, 91.02% %PF)

Run Status: Complete | Lane QC Status: QcPassed | Flow Cell Status: QcPassed

File Count/Size: 55,592 files (10 GB) | File Status: Active

Owner: PulseNet NGS Lab 1 | User: PulseNet NGS Lab 1

Latest Analysis: Cycles 151 | 10 | 10 | 151 | Yield 5.38 Gbp

- b. Click on the SampleSheet.csv link.

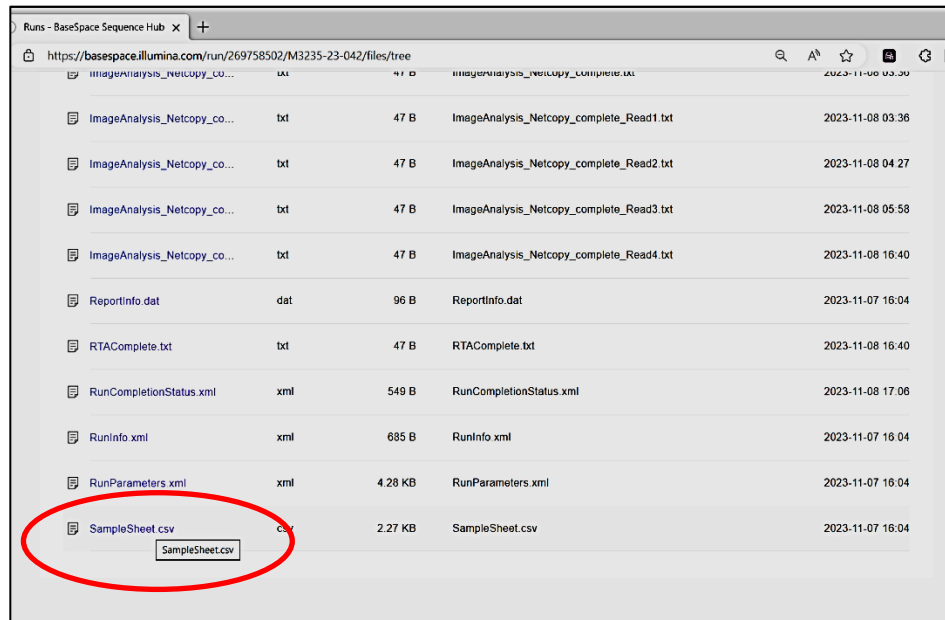
**PULSENET INTERNATIONAL STANDARD OPERATING PROCEDURE FOR ANALYZING ILLUMINA SHORT READ WGS DATA USING THE TERRA.BIO PLATFORM**

**Doc. No. PNID01**

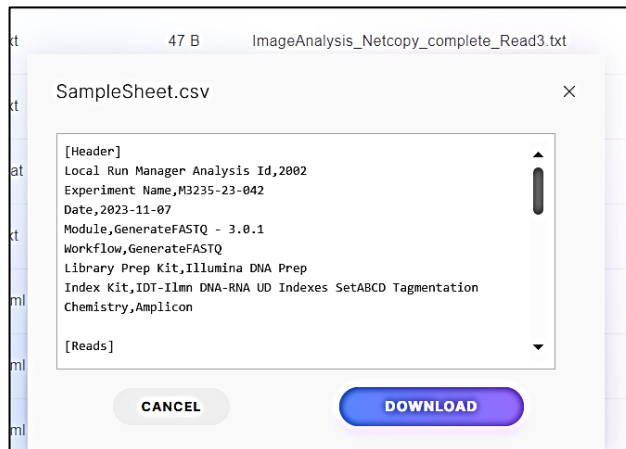
**Ver. No. 01**

**Effective Date:**

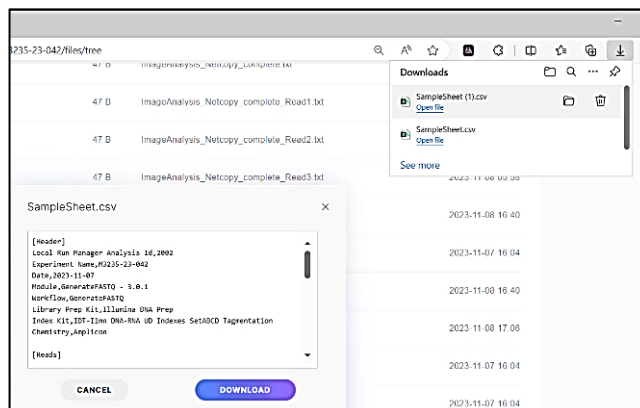
**Page 30 of 61**



c. In the “SampleSheet.csv” pop-up window, click on “Download”.



d. The downloaded csv file will appear on the top right corner under “Downloads”.



**PULSENET INTERNATIONAL STANDARD OPERATING PROCEDURE FOR ANALYZING ILLUMINA SHORT READ WGS DATA USING THE TERRA.BIO PLATFORM**

**Doc. No. PNID01**

**Ver. No. 01**

**Effective Date:**

**Page 31 of 61**

3. Open the SampleSheet. The columns that are needed in the metadata tsv file will depend on the columns and column content present in the SampleSheet.
4. Prepare the metadata tsv file:
  - a. Columns needed when the “Sample\_Name” and “Sample\_ID” columns in the SampleSheet have the exact **same** content:

Sample_ID	Sample_Name	Description	Index_Platform	Index	Index2	Sample_Project
D5480-M3235-23-042	D5480-M3235-23-042		B	A06	UDP0137	TATATTCGAG
ATCC-BAA-460-M3235-23-042	ATCC-BAA-460-M3235-23-042		B	B06	UDP0138	GGCCTTCGATA
2011L-2624-M3235-23-042	2011L-2624-M3235-23-042		B	C06	UDP0139	GAATACCT UDP0139
2015K-0092-M3235-23-042	2015K-0092-M3235-23-042		B	D06	UDP0140	TACGTGA UDP0140
2015K-1440-M3235-23-042	2015K-1440-M3235-23-042		B	E06	UDP0141	CTTATTGG UDP0141
2013V-1178-M3235-23-042	2013V-1178-M3235-23-042		B	F06	UDP0142	ACAACCTAC UDP0142
2014C-3598-M3235-23-042	2014C-3598-M3235-23-042		B	G06	UDP0143	GTTGGAT UDP0143
2014C-3857-M3235-23-042	2014C-3857-M3235-23-042		B	H06	UDP0144	AATCCAAT UDP0144
2015C-3881-M3235-23-042	2015C-3881-M3235-23-042		B	A07	UDP0145	TATGATG UDP0145

- i. **entity\_datatable\_name\_id**:
  1. Enter the name of the data table (either new or existing) into the cell A1 between “entity:” and “id”.
  2. Enter the sample IDs into the column A the way you want them to appear in the Terra data table.
- ii. **basespace\_sample\_name**: copy and paste the content from the SampleSheet “Sample\_Name” field.
- iii. **basespace\_collection\_id**: enter the Run name the way it appears on BaseSpace.

**PULSENET INTERNATIONAL STANDARD OPERATING PROCEDURE FOR ANALYZING ILLUMINA SHORT READ WGS DATA USING THE TERRA.BIO PLATFORM**

Doc. No. PNID01

Ver. No. 01

Effective Date:

Page 32 of 61

	A	B	C	D	E
1	entity:quality_control_id	basespace_sample_name	basespace_collection_id		
2	D5480-LRM4update	D5480-M3235-23-042	M3235-23-042		
3	ATCC-BAA-LRM4update	ATCC-BAA-460-M3235-23-042	M3235-23-042		
4	2011L-2624-LRM4update	2011L-2624-M3235-23-042	M3235-23-042		
5	2015K-0092-LRM4update	2015K-0092-M3235-23-042	M3235-23-042		
6	2015K-1440-LRM4update	2015K-1440-M3235-23-042	M3235-23-042		
7	2013V-1178-LRM4update	2013V-1178-M3235-23-042	M3235-23-042		
8	2014C-3598-LRM4update	2014C-3598-M3235-23-042	M3235-23-042		
9	2014C-3857-LRM4update	2014C-3857-M3235-23-042	M3235-23-042		
10	2015C-3881-LRM4update	2015C-3881-M3235-23-042	M3235-23-042		
11	2017C-3818-LRM4update	2017C-3818-M3235-23-042	M3235-23-042		
12	2017C-3830-LRM4update	2017C-3830-M3235-23-042	M3235-23-042		
13	2015C-5082-LRM4update	2015C-5082-M3235-23-042	M3235-23-042		

b. Columns needed when the “Sample\_Name” and “Sample\_ID” columns in the SampleSheet have **the different** content:

	A	B	C	D	E	F	G	H	I	J	K	L	M
1	[Header]												
2	Local Run Manager Analysis Id		123123										
3	Experiment Name	CAOC-M5870-230530E											
4	Date	6/14/2023											
5	Module	GenerateFASTQ_3.0.1											
6	Workflow	GenerateFASTQ											
7	Library Prep Kit	Illumina DNA Prep											
8	Index Kit	IDT-Illum DNA-RNA UD Indexes SetABCD Tagmentation											
9	Chemistry	Amplicon											
10													
11	[Reads]												
12		151											
13		151											
14													
15	[Settings]												
16	adapter	CTGTCTCTTATACATCT											
17													
18	[Data]												
19	Sample_ID	Sample_Name	Descripti	Index_Pla	Index_Pla	I7_Index	index	I5_Index	Index2	Sample_Project			
20	2023FD-00134	2023FD-00134-CAOC-M5870-230530E		A	A07	UDP0049	AGTGTGK/UDP0049	CTGGTAC/CPD_230530E					
21	BE230960535	BE230960535-CAOC-M5870-230530E		A	B07	UDP0050	GACACCA/UDP0050	TCACGTC/Sal_230530E					
22	2023FD-00135	2023FD-00135-CAOC-M5870-230530E		A	C07	UDP0051	CCTGTCTG/UDP0051	ACTGTGT/CPD_230530E					
23	2023FD-00136	2023FD-00136-CAOC-M5870-230530E		A	D07	UDP0052	TGATGTA/UDP0052	GTGGTGT/CPD_230530E					
24	2023FD-00137	2023FD-00137-CAOC-M5870-230530E		A	E07	UDP0053	GGAAITG/UDP0053	AGCACAT/CPD_230530E					
25	BE231320288	BE231320288-CAOC-M5870-230530E		A	F07	UDP0054	GCATAAG/UDP0054	TTCCGTG/ECOH_Shig230530E					
26	2023FD-00138	2023FD-00138-CAOC-M5870-230530E		A	G07	UDP0055	CTGAGGA/UDP0055	CTAACG/CPD_230530E					
27	2023FD-00139	2023FD-00139-CAOC-M5870-230530E		A	H07	UDP0056	AACGAC/UDP0056	GCCTCGG/CPD_230530E					
28	BE231330092	BE231330092-CAOC-M5870-230530E		A	A08	UDP0057	TCTATCT/UDP0057	CGTCGAC/Sal_230530E					
29	BE231330093	BE231330093-CAOC-M5870-230530E		A	B08	UDP0058	CTCGTTC/UDP0058	TACTAGT/Sal_230530E					
30	BE231350225	BE231350225-CAOC-M5870-230530E		A	C08	UDP0059	CTGTTGG/UDP0059	ATAGAC/Sal_230530E					

i. entity\_datatablename\_id:

1. Enter the name of the data table (either new or existing) into the cell A1 between “entity:” and “id”.
2. Enter the sample IDs into the column A the way you want them to appear in the Terra data table.

**PULSENET INTERNATIONAL STANDARD OPERATING PROCEDURE FOR ANALYZING ILLUMINA SHORT READ WGS DATA USING THE TERRA.BIO PLATFORM**

**Doc. No. PNID01**

**Ver. No. 01**

**Effective Date:**

**Page 33 of 61**

- ii. `basespace_sample_name`: copy and paste the content from the SampleSheet “Sample\_Name” field.
- iii. `basespace_sample_id`: copy and paste the content from the SampleSheet “Sample\_ID” field.
- iv. `basespace_collection_id`: enter the Run name the way it appears on BaseSpace.

	A	B	C	D	E	F	G	H	I
1	entity:quality_control_id	basespace_sample_name	basespace_sample_id	basespace_collection_id					
2	CAOC_2023FD-00134	2023FD-00134-CAOC-M5870-230530E	2023FD-00134	CAOC-M5870-230530E					
3	CAOC_BE230960535	BE230960535-CAOC-M5870-230530E	BE230960535	CAOC-M5870-230530E					
4	CAOC_2023FD-00135	2023FD-00135-CAOC-M5870-230530E	2023FD-00135	CAOC-M5870-230530E					
5	CAOC_2023FD-00136	2023FD-00136-CAOC-M5870-230530E	2023FD-00136	CAOC-M5870-230530E					
6	CAOC_2023FD-00137	2023FD-00137-CAOC-M5870-230530E	2023FD-00137	CAOC-M5870-230530E					
7									

**c. Columns needed for the NextSeq SampleSheet:**

	A	B	C
1	[Header]		
2	FileFormatVersion		2
3	RunName	VL403-24-001	
4	InstrumentPlatform	NextSeq1k2k	
5	IndexOrientation	Forward	
6			
7	[Reads]		
8	Read1Cycles		151
9	Read2Cycles		151
10	Index1Cycles		10
11	Index2Cycles		10
12			
13	[Sequencing_Settings]		
14	LibraryPrepKits	illuminaDNAPrep	
15			
16	[BCLConvert_Settings]		
17	SoftwareVersion	3.10.12	
18	AdapterRead1	CTGTCTCTTATACATCT	
19	AdapterRead2	CTGTCTCTTATACATCT	
20	OverrideCycles	Y151;I10;I10;Y151	
21	FastqCompressionFormat	gzip	
22			
23	[BCLConvert_Data]		
24	Sample_ID	Index	Index2
25	2013L-5214	CGACATCCGA	TACGTTTCATT
26	2013L-5351	GCACAATAGGA	TCCATCCGAG
27	2013L-5356	GCACAATAGGA	CTTGTCTTAA
28	2013L-5357	TAGTTCCGTA	CCATGTGTAG
29	2013L-5585	CTATTACTAC	GAGTCTCTCC
30	2013L-5615	TAGCATAACC	GCTATGGCCA
31	2011L-2624	ACTCTATTGT	ATCGCATATG
32	2015K-9887	CCAAAGGCTT	TGGAAGTACT
33	2015K-1104	TACTCCACA	GACACCGATG
34	2014K-0833	AGTAGAAGTG	CTAGCGTCCA

- i. `entity_datatable_name_id`:
  - 1. Enter the name of the data table (either new or existing) into the cell A1 between “entity:” and “id”.
  - 2. Enter the sample IDs into the column A the way you want them to appear in the Terra data table.
- ii. `basespace_sample_name`: copy and paste the content from the SampleSheet “Sample\_ID” field.
- iii. `basespace_sample_id`: copy and paste the content from the SampleSheet “Sample\_ID” field.

**PULSENET INTERNATIONAL STANDARD OPERATING PROCEDURE FOR ANALYZING ILLUMINA SHORT READ WGS DATA USING THE TERRA.BIO PLATFORM**

**Doc. No. PNID01**

**Ver. No. 01**

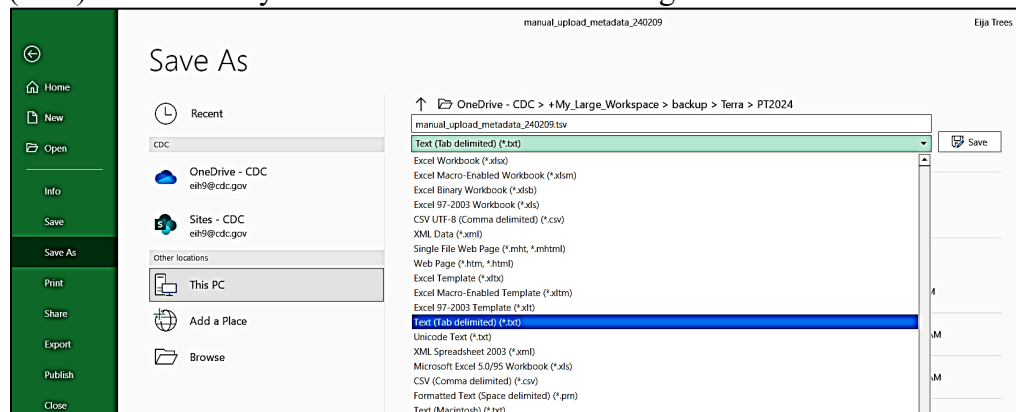
**Effective Date:**

**Page 34 of 61**

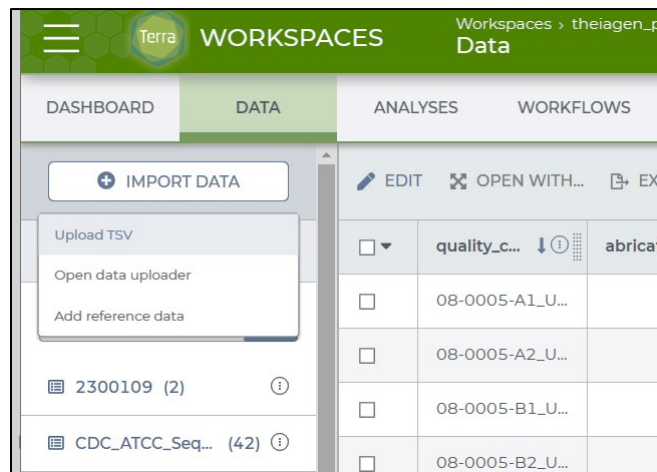
- iv. `basespace_collection_id`: enter the Run name the way it appears on BaseSpace.

entity:quality_control_id	basespace_sample_name	basespace_sample_id	basespace_collection_id
2013L-5214_v3index_NextSeq	2013L-5214	2013L-5214	VL403-24-001
2013L-5351_v3index_NextSeq	2013L-5351	2013L-5351	VL403-24-001
2013L-5356_v3index_NextSeq	2013L-5356	2013L-5356	VL403-24-001
2013L-5357_v3index_NextSeq	2013L-5357	2013L-5357	VL403-24-001
2013L-5585_v3index_NextSeq	2013L-5585	2013L-5585	VL403-24-001
2013L-5615_v3index_NextSeq	2013L-5615	2013L-5615	VL403-24-001
2011L-2624_v3index_NextSeq	2011L-2624	2011L-2624	VL403-24-001
2015K-0887_v3index_NextSeq	2015K-0887	2015K-0887	VL403-24-001
2015K-1104_v3index_NextSeq	2015K-1104	2015K-1104	VL403-24-001
2014K-0833_v3index_NextSeq	2014K-0833	2014K-0833	VL403-24-001

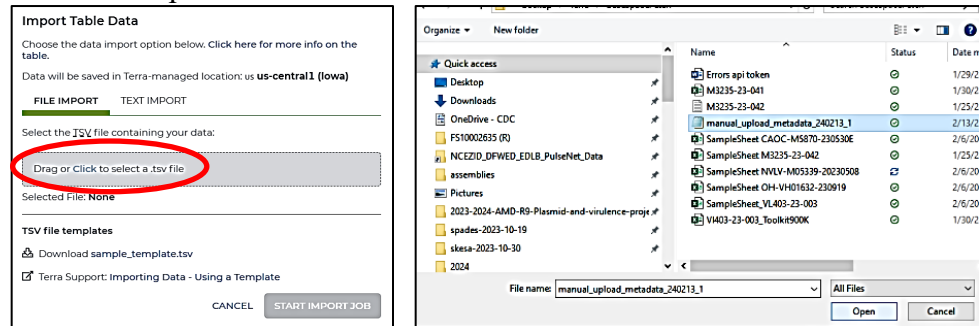
- d. Save the file in the **tsv format**: choose “Save As” and “Text (Tab delimited) (\*.txt)”. Make sure your file name has “.tsv” ending.



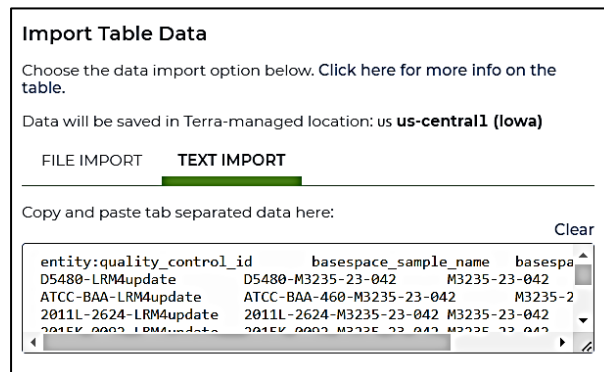
- 5. Import the metadata tsv file:
  - a. In the “Data” tab, click on “Import Data” and select “Upload tsv” from the drop-down menu.



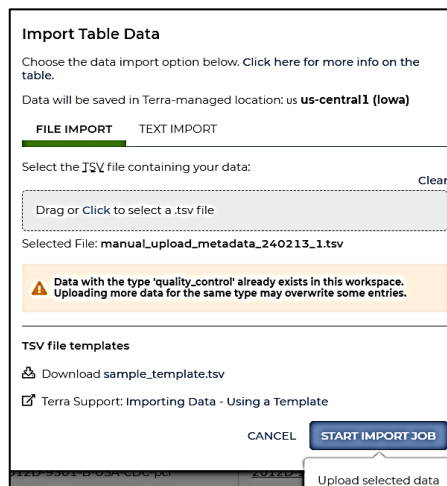
- b. From the “Import Table Data” pop-up window, you can import the metadata in two ways:
  - i. In the “File Import” tab, click in the middle to select the tsv file and navigate to the location where the metadata tsv file is saved, select the file and click “Open”.



- ii. Alternatively, you can switch to the “text Import” tab and copy and paste the contents of the tsv file into the field in the middle.



- c. In the “Import Table Data” pop-up window, you will get a warning that data already exists in the data table in question (if importing to an existing data table) and uploading more data to it may overwrite existing data. Click on “Start Import Job”.



**PULSENET INTERNATIONAL STANDARD OPERATING PROCEDURE FOR ANALYZING ILLUMINA SHORT READ WGS DATA USING THE TERRA.BIO PLATFORM**

**Doc. No. PNID01**

**Ver. No. 01**

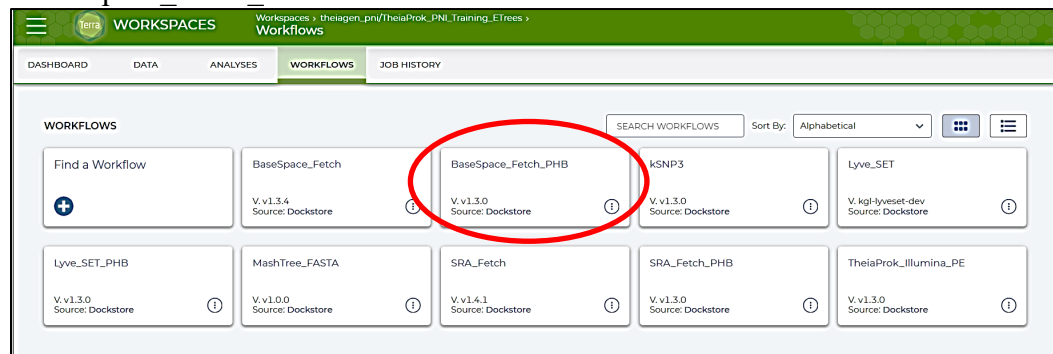
**Effective Date:**

**Page 36 of 61**

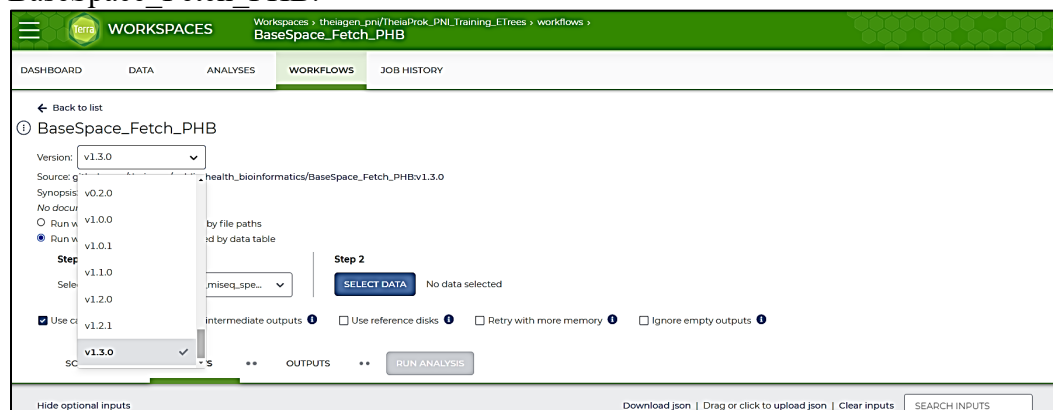
- d. After the import is done you should see new entries created to the data table for the sequences to be imported from BaseSpace.

quality_control_id	basespace_collection_id	basespace_fetch_analysis_date	basespace_fetch_vers
D5480-LRM4update	M3235-23-042	2024-01-30	Terra Utilities v1.3.4
D7320-L-1A_USA_CDC_pcl			
D7320-L-1B_USA_CDC_pcl			
D7320-L-2A_USA_CDC_pcl			
D7320-L-2B_USA_CDC_pcl			
NVLV_QA-157	NVLV-M05339-20230508		
NVLV_QA-158	NVLV-M05339-20230508		
NVLV_QA-159	NVLV-M05339-20230508		
NVLV_QA-160	NVLV-M05339-20230508		
NVLV_QA-161	NVLV-M05339-20230508		

6. Run the “BaseSpace\_Fetch\_PHB” workflow:
  - a. In the “Workflows” tab, click on “BaseSpace\_Fetch\_PHB”. This will open the “BaseSpace\_Fetch\_PHB” screen.



- b. From the “Version” drop-down menu, select the latest version of BaseSpace\_Fetch\_PHB.



- c. Under “Step 1”, click on the “Select root entity type” drop-down menu and select the data table into which you uploaded (steps 4-5) the tsv file containing the metadata for the run to be imported from BaseSpace.

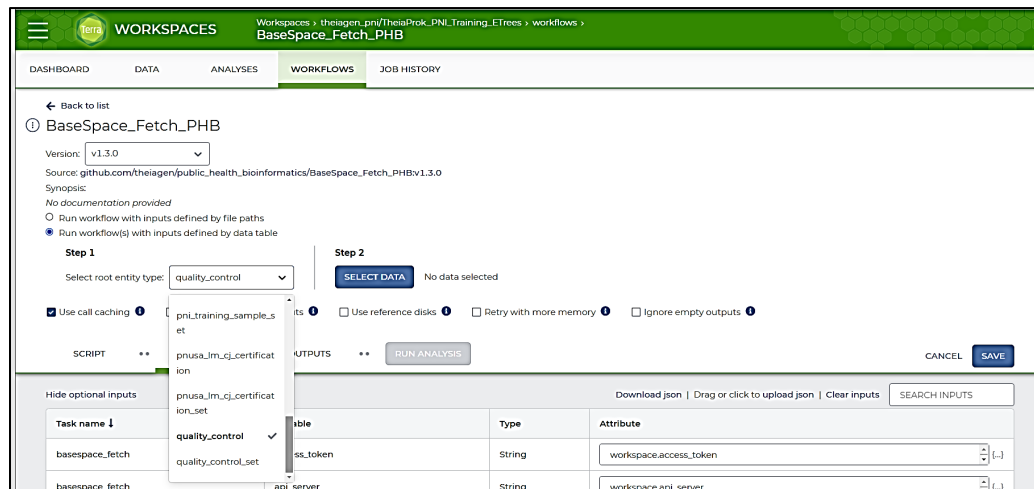
**PULSENET INTERNATIONAL STANDARD OPERATING PROCEDURE FOR ANALYZING ILLUMINA SHORT READ WGS DATA USING THE TERRA.BIO PLATFORM**

**Doc. No. PNID01**

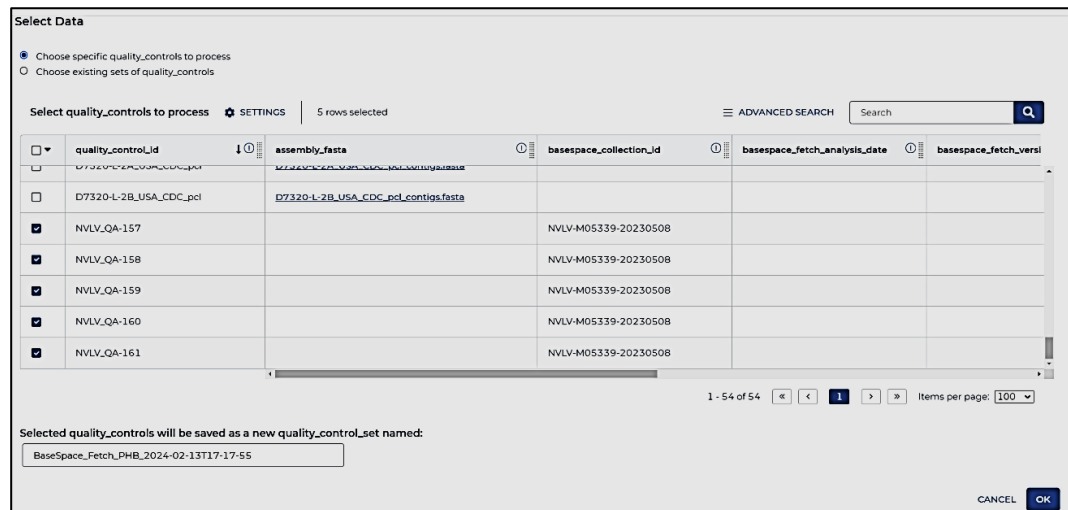
**Ver. No. 01**

**Effective Date:**

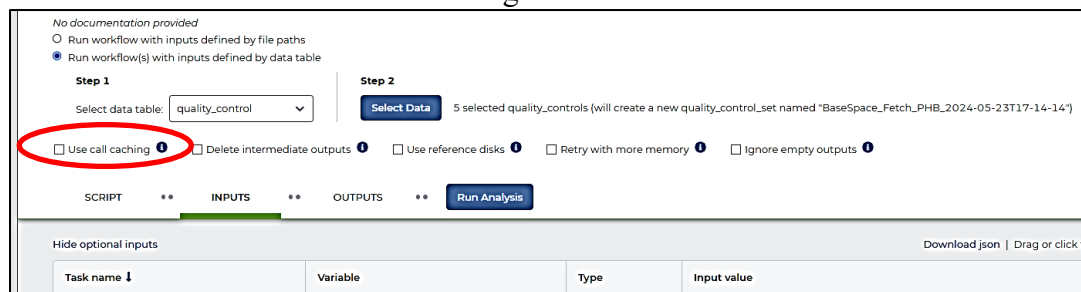
**Page 37 of 61**



- d. Under “Step 2”, click on “Select data” (screenshot above). This will take you to the sample selection screen.
- e. Check the boxes next to the samples to be imported from BaseSpace and click “OK”. This will take you back to the “BaseSpace\_Fetch\_PHB” screen.



- f. De-select the box for “Use call catching”.



- g. In the “Inputs” tab, define the following variables in the “Attribute” field:  
**NOTE:** When you fill in the Attribute column, clicking inside the cell will bring up a drop-down menu of attributes that you can select to avoid typos
  - i. access\_token: “workspace.access\_token”.

**PULSENET INTERNATIONAL STANDARD OPERATING PROCEDURE FOR ANALYZING ILLUMINA SHORT READ WGS DATA USING THE TERRA.BIO PLATFORM**

**Doc. No. PNID01**

**Ver. No. 01**

**Effective Date:**

**Page 38 of 61**

- ii. api\_server: “workspace.api\_server”.
- iii. basespace\_collection\_id: ”this.basespace\_collection\_id”.
- iv. basespace\_sample\_name: “this.basespace\_sample\_name”.
- v. sample\_name: “**this.datatablename\_id**”, e.g. “this.quality\_control\_id”.
- vi. basespace\_sample\_id: “this.basespace\_sample\_id”.

**NOTE:** *needs to be filled out **only** if the contents in the “Sample\_Name” and “Sample\_ID” fields are different in the run SampleSheet or you are importing data from a NextSeq run.*

Task name ↓	Variable	Type	Attribute
basespace_fetch	access_token	String	workspace.access_token
basespace_fetch	api_server	String	workspace.api_server
basespace_fetch	basespace_collection_id	String	this.basespace_collection_id
basespace_fetch	basespace_sample_name	String	this.basespace_sample_name
basespace_fetch	sample_name	String	this.quality_control_id
basespace_fetch	basespace_sample_id	String	this.basespace_sample_id
fetch_bs	cpu	Int	Optional
fetch_bs	disk_size	Int	Optional

- h. In the “**Outputs**” tab, click “Use defaults”, then click “Save”, and then “Run Analysis”.

**NOTE:** *The “Save” button is only visible if parameters (other than sample IDs) have changed from the previous job submission. The “Run Analysis” button becomes highlighted after you save the parameters.*

Output files will be saved to  
 Files / submission unique ID / basespace\_fetch / workflow unique ID

References to outputs will be written to  
 Tables / quality\_control

Fill in the attributes below to add or update columns in your data table

Task name ↓	Variable	Type	Attribute   Use defaults ←
basespace_fetch	basespace_fetch_analysis_date	String	this.basespace_fetch_analysis_date
basespace_fetch	basespace_fetch_version	String	this.basespace_fetch_version
basespace_fetch	read1	File	this.read1
basespace_fetch	read2	File	this.read2

- i. In the “Confirm launch” pop-up window, describe your job submission (optional) and click “Launch”.

**PULSENET INTERNATIONAL STANDARD OPERATING PROCEDURE FOR ANALYZING ILLUMINA SHORT READ WGS DATA USING THE TERRA.BIO PLATFORM**

Doc. No. PNID01

Ver. No. 01

Effective Date:

Page 39 of 61

**Confirm launch**

Output files will be saved as workspace data in:  
us us-central1 (Iowa) ⓘ

Running workflows will generate cloud charges. ⓘ  
How much does my workflow cost? ⓘ  
Set up budget alert ⓘ

Describe your submission (optional):

MiSeq run M3235-23-042 from BaseSpace

This will launch **20** analyses.

CANCEL LAUNCH

- j. The “Job History” tab opens where the status of the job submissions should initially be “Queued”.

WORKSPACES  
Job History

← Job History › Submission 497ec4b7-28b0-47d6-9140-903cd92305cb

<b>Workflow Statuses</b> Submitted: 5	<b>Workflow Configuration</b> theiaigen_pni/BaseSpace_Fetch_PHB	<b>Submitted by</b> eja.trees@theiaigen.cloud Feb 13, 2024, 12:25 PM	<b>Total Run Cost</b> N/A
<b>Data Entity</b> BaseSpace_Fetch_PHB_2024-02-13T17-17-55 quality_control_set	<b>Submission ID</b> 497ec4b7-28b0-47d6-9140-903cd92305cb	<b>Call Caching</b> Enabled	
<b>Comment</b> NVLV MiSeq Basespace run. Sample_na...	<b>Delete Intermediate Outputs</b> Disabled	<b>Use Reference Disks</b> Disabled	<b>Retry with More Memory</b> Disabled

**WORKFLOWS** | INPUTS | OUTPUTS

Search workflows:  Completion status:  Download TSV

Data Entity ↓	Last Changed	Status	Run Cost	Messages	Workflow ID
NVLV_QA-157 (quality_control)	Feb 13, 2024, 12:25 PM	⌚ Queued	N/A		
NVLV_QA-158 (quality_control)	Feb 13, 2024, 12:25 PM	⌚ Queued	N/A		
NVLV_QA-159 (quality_control)	Feb 13, 2024, 12:25 PM	⌚ Queued	N/A		
NVLV_QA-160 (quality_control)	Feb 13, 2024, 12:25 PM	⌚ Queued	N/A		
NVLV_QA-161 (quality_control)	Feb 13, 2024, 12:25 PM	⌚ Queued	N/A		

- k. The status will be “Succeeded” or “Done” once the job has finished. In the “Data” tab, you should now see the FASTQ file names in “Read1” and “Read2” columns and also information in columns “basespace\_fetch\_analysis\_date” and “basespace\_fetch\_version”.

# PULSENET INTERNATIONAL STANDARD OPERATING PROCEDURE FOR ANALYZING ILLUMINA SHORT READ WGS DATA USING THE TERRA.BIO PLATFORM

**Doc. No. PNID01**

**Ver. No. 01**

**Effective Date:**

**Page 40 of 61**

← Job History > Submission 497ec4b7-28b0-47d6-9140-903cd92305cb

**Workflow Statuses**

✓ Succeeded: 5

**Workflow Configuration**

theigen\_pnl/Basespace\_Fetch\_PHB

**Submitted by**

eija.trees@theigencloud  
Feb 13, 2024, 12:25 PM

**Total Run Cost**

N/A

**Data Entity**

Basespace\_Fetch\_PHB.2024-02-13T17:17:55  
quality\_control.set

**Submission ID**

497ec4b7-28b0-47d6-9140-903cd92305cb

**Call Caching**

Enabled

**Use Reference Disks**

Disabled

**Comment**

NVLV\_MISEq Basespace run. Sample\_na...

**Delete Intermediate Outputs**

Disabled

**Retry with More Memory**

Disabled

**WORKFLOWS**    INPUTS    OUTPUTS

Completion status
Download TSV

Data Entity ↓	Last Changed	Status	Run Cost	Messages	Workflow ID	Links
NVLV_QA-157 (quality_control)	Feb 13, 2024, 12:30 PM	✓ Succeeded	N/A		6154609d-c6dc-4d1d-e1e8-5286b71303...	<a href="#">🔗</a> <a href="#">📄</a> <a href="#">📁</a>
NVLV_QA-158 (quality_control)	Feb 13, 2024, 12:29 PM	✓ Succeeded	N/A		d3535af4-758a-490d-9be9-16696457800...	<a href="#">🔗</a> <a href="#">📄</a> <a href="#">📁</a>
NVLV_QA-159 (quality_control)	Feb 13, 2024, 12:29 PM	✓ Succeeded	N/A		1e1351e2-582d-4978-9f77-0ecba6e59472...	<a href="#">🔗</a> <a href="#">📄</a> <a href="#">📁</a>
NVLV_QA-160 (quality_control)	Feb 13, 2024, 12:29 PM	✓ Succeeded	N/A		f56171f1f-Gada-455c-8820-8f6fb54b4094...	<a href="#">🔗</a> <a href="#">📄</a> <a href="#">📁</a>
NVLV_QA-161 (quality_control)	Feb 13, 2024, 12:29 PM	✓ Succeeded	N/A		3be0f9ce-a2e3-4302-af5a-01187ac5eb07...	<a href="#">🔗</a> <a href="#">📄</a> <a href="#">📁</a>

Workspaces > theigen\_pnl/TheiaProk\_PNI\_Training\_ETrees > Data

DASHBOARD    DATA    ANALYSES    WORKFLOWS    JOB HISTORY

illumina\_pe\_v1... (25)

analysis\_pt\_24 (30)

analysis\_pt\_24.set (3)

orange\_miseq\_s... (20)

orange\_miseq.sp... (1)

pnl\_training\_sa... (13)

pnl\_training.sam... (1)

pnusa\_lm\_cj.cer... (96)

pnusa\_lm\_cj.cer... (29)

**quality\_control (54)**

quality\_control... (12)

EDIT    OPEN WITH...    EXPORT    SETTINGS    0 rows selected    ADVANCED SEARCH    Search

quality_controlId	e_sample_name	read1	read2
D5480-LRM4update	I3235-23-042	D5480-LRM4update_R1.fastq.gz	D5480-LRM4update_R2.fastq.gz
D7320-L-1A_USA_CDC.pcd		D7320-L-1A-M3235-21-007-512_L001_R1_001.fast...	D7320-L-1A-M3235-21-007-512_L001_...
D7320-L-1B_USA_CDC.pcd		D7320-L-1B-M3235-21-007-513_L001_R1_001.fast...	D7320-L-1B-M3235-21-007-513_L001_...
D7320-L-2A_USA_CDC.pcd		D7320-L-2A-M3235-21-007-514_L001_R1_001.fast...	D7320-L-2A-M3235-21-007-514_L001_...
D7320-L-2B_USA_CDC.pcd		D7320-L-2B-M3235-21-007-515_L001_R1_001.fast...	D7320-L-2B-M3235-21-007-515_L001_...
NVLV_QA-157		NVLV_QA-157_R1.fastq.gz	NVLV_QA-157_R2.fastq.gz
NVLV_QA-158		NVLV_QA-158_R1.fastq.gz	NVLV_QA-158_R2.fastq.gz
NVLV_QA-159		NVLV_QA-159_R1.fastq.gz	NVLV_QA-159_R2.fastq.gz
NVLV_QA-160		NVLV_QA-160_R1.fastq.gz	NVLV_QA-160_R2.fastq.gz
NVLV_QA-161		NVLV_QA-161_R1.fastq.gz	NVLV_QA-161_R2.fastq.gz

1 - 54 of 54    Items per page: 100

Workspaces > theigen\_pnl/TheiaProk\_PNI\_Training\_ETrees > Data

DASHBOARD    DATA    ANALYSES    WORKFLOWS    JOB HISTORY

illumina\_pe\_v1... (25)

analysis\_pt\_24 (30)

analysis\_pt\_24.set (3)

orange\_miseq\_s... (20)

orange\_miseq.sp... (1)

pnl\_training\_sa... (13)

pnl\_training.sam... (1)

pnusa\_lm\_cj.cer... (96)

pnusa\_lm\_cj.cer... (29)

**quality\_control (54)**

quality\_control... (12)

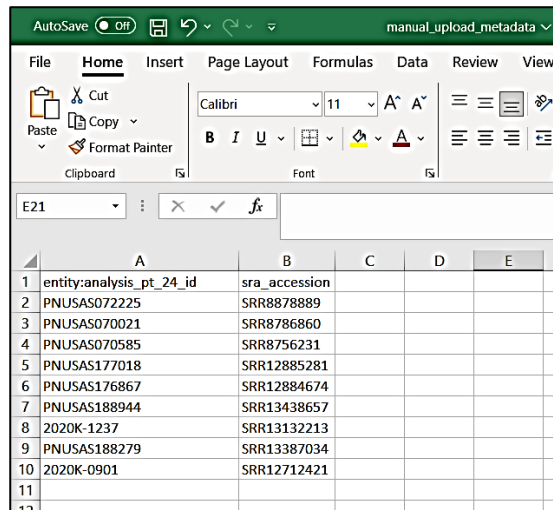
EDIT    OPEN WITH...    EXPORT    SETTINGS    0 rows selected    ADVANCED SEARCH    Search

quality_controlId	basespace_fetch_analysis_date	basespace_fetch_version	basespace_sample_name
D5480-LRM4update	2024-01-30	Terra Utilities v1.3.4	D5480-M3235-23-042
D7320-L-1A_USA_CDC.pcd			
D7320-L-1B_USA_CDC.pcd			
D7320-L-2A_USA_CDC.pcd			
D7320-L-2B_USA_CDC.pcd			
NVLV_QA-157	2024-02-13	PHB v1.3.0	QA-157
NVLV_QA-158	2024-02-13	PHB v1.3.0	QA-158
NVLV_QA-159	2024-02-13	PHB v1.3.0	QA-159
NVLV_QA-160	2024-02-13	PHB v1.3.0	QA-160
NVLV_QA-161	2024-02-13	PHB v1.3.0	QA-161

1 - 54 of 54    Items per page: 100

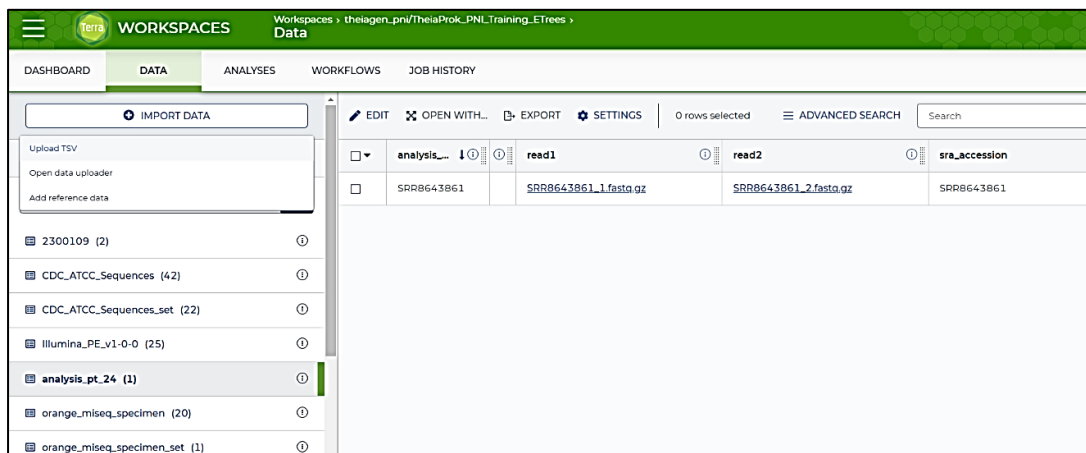
**Appendix PNID01-2: Data Download from NCBI SRA**

1. Prepare the metadata tsv file:
  - a. Enter the name of the data table (either new or existing) into the cell A1 between “entity:” and “id”.
  - b. Enter the sample IDs into the column A the way you want them to appear in the Terra data table.
  - c. sra\_accession: enter the SRA accession numbers for the sequences to be downloaded from SRA.

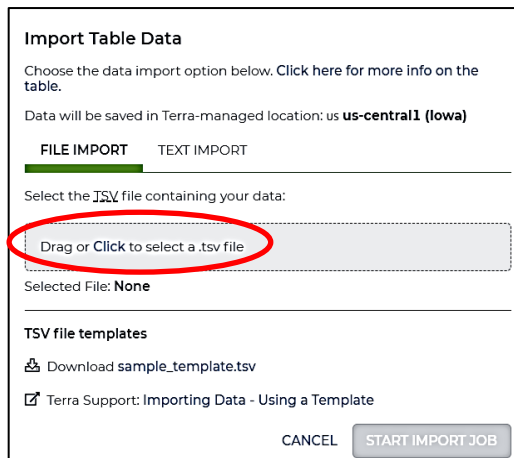


- d. Save the file in the **tsv format**: choose “Save As” and “Text (Tab delimited) (\*.txt)”. Make sure your file name has “.tsv” ending.

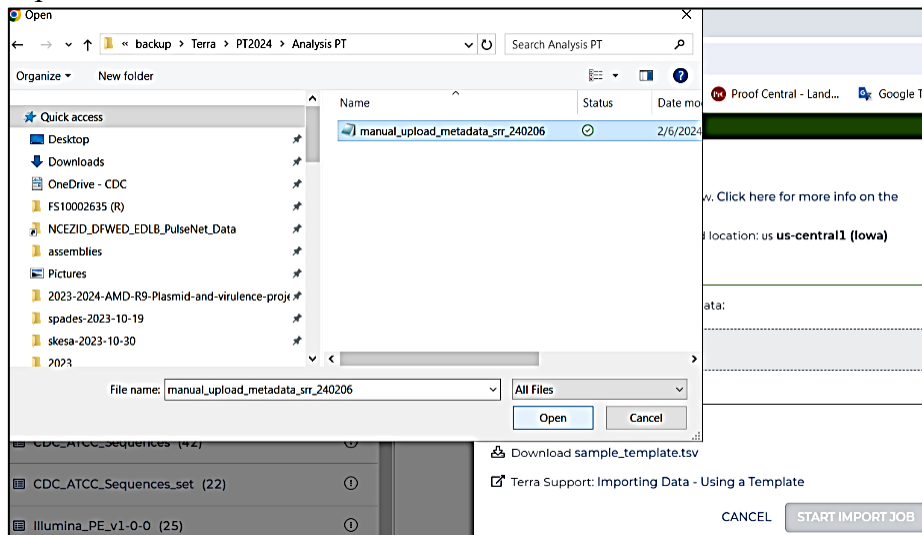
2. In the “Data” tab, click on “Import Data” and select “Upload tsv” from the drop-down menu.



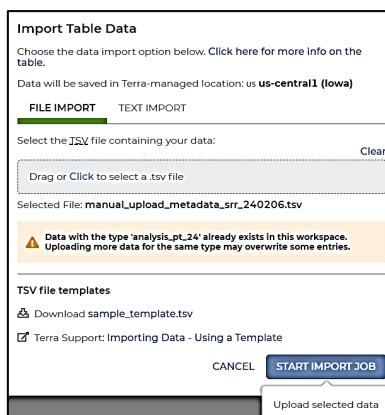
3. In the “Import Table Data” pop-up window “File Import” tab, click in the middle to select the tsv file.



4. Navigate to the location where the metadata tsv file is saved, select the file and click “Open”.



5. In the “Import Table Data” pop-up window, you will get a warning that data already exists in the data table in question (if importing to an existing data table) and uploading more data to it may overwrite existing data. Click on “Start Import Job”.



**PULSENET INTERNATIONAL STANDARD OPERATING PROCEDURE FOR ANALYZING ILLUMINA SHORT READ WGS DATA USING THE TERRA.BIO PLATFORM**

**Doc. No. PNID01**

**Ver. No. 01**

**Effective Date:**

**Page 43 of 61**

- After the import is done you should see new entries created to the data table for the sequences to be downloaded from NCBI together with their SRA accession numbers

analysis...	read1	read2	sra_accession
2020k-0901			SRR12712421
2020k-1237			SRR13132213
PNUSAS0700...			SRR8786860
PNUSAS0705...			SRR8756231
PNUSAS0722...			SRR8878889
PNUSAS1768...			SRR12884674
PNUSAS1770...			SRR12885281
PNUSAS1882...			SRR13387034
PNUSAS1889...			SRR13438657

- In the “Workflows” tab, click on the “SRA\_Fetch\_PHB” workflow. This will open the “SRA\_Fetch\_PHB” screen.

The screenshot shows the 'WORKFLOWS' tab in Terra Bio Workspaces. A grid of workflow cards is displayed, including 'Find a Workflow', 'BaseSpace\_Fetch', 'BaseSpace\_Fetch\_PHB', 'kSNP3', 'Lyve\_SET', 'Lyve\_SET\_PHB', 'MashTree\_FASTA', 'SRA\_Fetch', 'SRA\_Fetch\_PHB' (circled in red), and 'TheiaProk\_Illumina\_PE'. Each card shows the workflow name, version, and source (Dockstore).

- From the “Version” drop-down menu, select the latest version of SRA\_Fetch\_PHB.

The screenshot shows the configuration screen for the 'SRA\_Fetch\_PHB' workflow. The 'Version' dropdown menu is open, displaying a list of versions: v0.2.0, v1.0.0, v1.0.1, v1.1.0, v1.2.0, v1.2.1, and v1.3.0. The v1.3.0 version is selected. The 'Step 2' section shows a 'SELECT DATA' button and a 'RUN ANALYSIS' button. The 'Use reference disks' checkbox is checked.

**PULSENET INTERNATIONAL STANDARD OPERATING PROCEDURE FOR ANALYZING ILLUMINA SHORT READ WGS DATA USING THE TERRA.BIO PLATFORM**

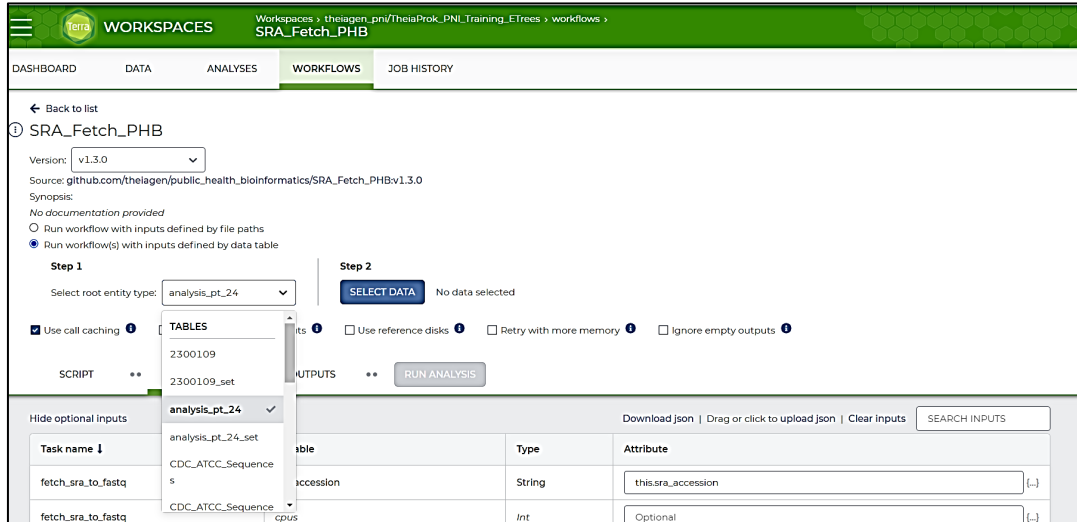
**Doc. No. PNID01**

**Ver. No. 01**

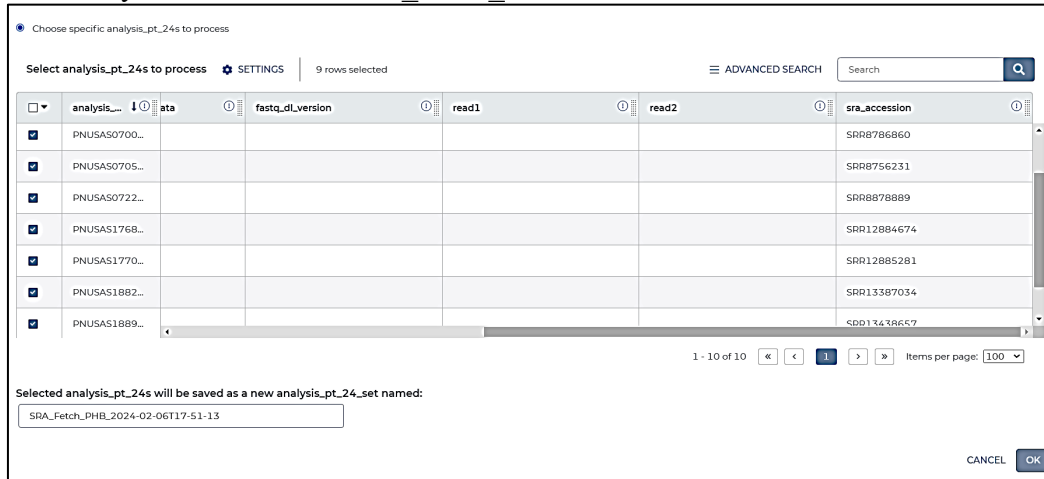
**Effective Date:**

**Page 44 of 61**

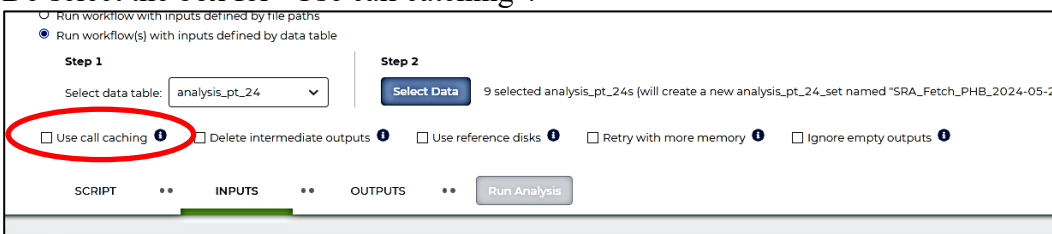
- Under “Step 1”, click on the “Select root entity type” drop-down menu and select the data table into which you imported (steps 1-6) the tsv file containing the SRA accession numbers for the samples to be downloaded from SRA.



- Under “Step 2”, click on “Select data” (screenshot above). This will take you to the sample selection screen.
- Check the boxes next to the samples to be downloaded from NCBI and click “OK”. This will take you back to the “SRA\_Fetch\_PHB” screen.



- De-select the box for “Use call caching”.



- In the “Inputs” tab, define the “sra\_accession” variable in the “Attribute” field: “this.sra\_accession”.

**PULSENET INTERNATIONAL STANDARD OPERATING PROCEDURE FOR ANALYZING ILLUMINA SHORT READ WGS DATA USING THE TERRA.BIO PLATFORM**

**Doc. No. PNID01**

**Ver. No. 01**

**Effective Date:**

**Page 45 of 61**

Task name ↓	Variable	Type	Attribute
fetch_sra_to_fastq	sra_accession	String	<input type="text" value="this.sra_accession"/> [-]
fetch_sra_to_fastq	cpus	Int	<input type="text" value="Optional"/> [-]
fetch_sra_to_fastq	disk_size	Int	<input type="text" value="Optional"/> [-]

14. In the “**Outputs**” tab, click “Use defaults”, then click “Save”, and then “Run Analysis”.  
**NOTE:** The “Save” button is only visible if parameters (other than sample IDs) have changed from the previous job submission.

Task name ↓	Variable	Type	Attribute
fetch_sra_to_fastq	fastq_dl_date	String	<input type="text" value="this.fastq_dl_date"/> [-]
fetch_sra_to_fastq	fastq_dl_docker	String	<input type="text" value="this.fastq_dl_docker"/> [-]
fetch_sra_to_fastq	fastq_dl_fastq_metadata	File	<input type="text" value="this.fastq_dl_fastq_metadata"/> [-]

15. In the “**Confirm launch**” pop-up window, describe your job submission (optional) and click “Launch”.

**Confirm launch**

Output files will be saved as workspace data in:  
us-central1 (Iowa) ⓘ

Running workflows will generate cloud charges. ⓘ  
How much does my workflow cost? ⌵  
Set up budget alert ⌵

Describe your submission (optional):

This will launch 9 analyses.

**CANCEL LAUNCH**

16. The “**Job History**” tab opens where the status of the job submissions should initially be “**Queued**”.

# PULSENET INTERNATIONAL STANDARD OPERATING PROCEDURE FOR ANALYZING ILLUMINA SHORT READ WGS DATA USING THE TERRA.BIO PLATFORM

Doc. No. PNID01

Ver. No. 01

Effective Date:

Page 46 of 61

The screenshot shows the 'Job History' page in Terra Bio Workspaces. The top navigation bar includes 'DASHBOARD', 'DATA', 'ANALYSES', 'WORKFLOWS', and 'JOB HISTORY'. The main content area displays workflow details for submission ID 86b637b7-edfb-4fbe-8b9e-d523773fadc. Key information includes: Workflow Status (Submitted: 9), Workflow Configuration (theiagen\_pni/SRA\_Fetch\_PHB), Submitted by (eja.trees@theiagen.cloud), Total Run Cost (N/A), Data Entity (SRA\_Fetch\_PHB\_2024-02-06T17-51-13 analysis\_pt\_24\_set), Submission ID (86b637b7-edfb-4fbe-8b9e-d523773fadc), Call Caching (Enabled), Comment (Analysis PT candidate set 1 (Newport)), Delete Intermediate Outputs (Disabled), Use Reference Disks (Disabled), and Retry with More Memory (Disabled). Below this, there is a 'WORKFLOWS' section with a search bar and a table listing workflow runs.

Data Entity	Last Changed	Status	Run Cost	Messages	Workflow ID	Links
2020K-0901 (analysis_pt_24)	Feb 6, 2024, 12:56 PM	⏸ Queued	N/A			
2020K-1237 (analysis_pt_24)	Feb 6, 2024, 12:56 PM	⏸ Queued	N/A			
PNUSAS070021 (analysis_pt_24)	Feb 6, 2024, 12:56 PM	⏸ Queued	N/A			

17. The status will be “Done” once the job has finished. In the “Data” tab, you should now see the FASTQ file names in “Read1” and “Read2” columns.

This screenshot shows the 'Job History' page with a table listing workflow runs. The status for the selected runs is 'Done'. A tooltip is visible over the 'SRA\_Fetch\_PHB\_2024-02-06T17-51-13 (analysis\_pt\_24\_set)' entry.

Submission (click for details)	Data entity	No. of Workflows	Status	Submitted	Submission ID	Comment	Actions
SRA_Fetch_PHB Submitted by eja.trees@theiagen.cloud	SRA_Fetch_PHB_2024-0...	9	✓ Done	Feb 6, 2024 12:56 PM	86b637b7-edfb-4fbe-8b9e-d523773fadc	Analysis PT candidate set ...	ⓘ
BaseSpace_Fetch_PHB Submitted by curtis.kapsak@theiagen.com	SRA_Fetch_PHB_2024-02-06T17-51-13 (analysis_pt_24_set)		✓ Done	Feb 2, 2024 3:15 PM	079d5e8e-100b-4977-a6b5-4de5baf9a8aa	test on OC miseq run wh...	ⓘ

This screenshot shows the 'Data' page in Terra Bio Workspaces. It displays a table of FASTQ files. The table has columns for 'analysis\_pt\_24...', 'read1', 'read2', and 'sra\_accession'. The 'read1' and 'read2' columns contain FASTQ file names (e.g., SRR12712421\_1.fastq.gz). The 'sra\_accession' column contains SRA accession numbers (e.g., SRR12712421). A left sidebar shows a list of tables, with 'analysis\_pt\_24 (10)' selected.

analysis_pt_24...	read1	read2	sra_accession
2020K-0901	SRR12712421_1.fastq.gz	SRR12712421_2.fastq.gz	SRR12712421
2020K-1237	SRR13132213_1.fastq.gz	SRR13132213_2.fastq.gz	SRR13132213
PNUSAS068804	SRR8643861_1.fastq.gz	SRR8643861_2.fastq.gz	SRR8643861
PNUSAS070021	SRR8786860_1.fastq.gz	SRR8786860_2.fastq.gz	SRR8786860
PNUSAS070585	SRR8756231_1.fastq.gz	SRR8756231_2.fastq.gz	SRR8756231
PNUSAS072225	SRR8878889_1.fastq.gz	SRR8878889_2.fastq.gz	SRR8878889
PNUSAS176867	SRR12884674_1.fastq.gz	SRR12884674_2.fastq.gz	SRR12884674
PNUSAS177018	SRR12885281_1.fastq.gz	SRR12885281_2.fastq.gz	SRR12885281
PNUSAS188279	SRR13387034_1.fastq.gz	SRR13387034_2.fastq.gz	SRR13387034
PNUSAS188944	SRR13438657_1.fastq.gz	SRR13438657_2.fastq.gz	SRR13438657

**Appendix PNID01-3: Customization of a Data Table View for PulseNet QC Metrics**

**NOTE:** *this customization can only be done after you have run the TheiaProk workflow once.*

1. In the “Data” tab, select the data table of interest, e.g., “CDC\_ATCC\_Sequences”, then select “Settings”.
2. Under “Select columns” the following QC metrics should be checked:
  - a. Ani\_highest\_percent
  - b. Ani\_top\_species\_match
  - c. Assembly\_length
  - d. Combined\_mean\_q\_clean
  - e. Combined\_mean\_q\_raw
  - f. Combined\_mean\_readlength\_clean
  - g. Combined\_mean\_readlength\_raw
  - h. Est\_coverage\_clean
  - i. Est\_coverage\_raw
  - j. Gambit\_predicted\_taxon
  - k. Midas\_secondary\_genus
  - l. Midas\_secondary\_genus\_abundance
  - m. N50\_value
  - n. Number\_contigs
  - o. Raw\_read\_screen
  - p. Seqsero2\_predicted\_contamination

**Select columns**

Show: all | none      Sort: alphabetical

- agrvate\_summary
- agrvate\_version
- combined\_mean\_q\_clean
- combined\_mean\_q\_raw
- combined\_mean\_readlength\_clean
- combined\_mean\_readlength\_raw
- meningotype\_BAST
- meningotype\_FetA
- meningotype\_NHBA
- meningotype\_NadA
- meningotype\_PorA
- meningotype\_PorB
- meningotype\_fHbp

SAVE THIS COLUMN SELECTION

Your saved column selections:

- pulsenet\_genotyping ⓘ
- qc\_metrics ⓘ

CANCEL      **DONE**

3. Click “Save this column selection”, name the column selection “qc\_metrics” and click “Save” and then click “Done”.

**NOTE:** *If you are adding or deleting columns from an existing column selection, click “Save this column selection”, select the name from the drop-down menu and click “Update”.*

**Creating a new column selection**

The screenshot shows the 'Select columns' dialog box with the following elements:

- Show:** all | none
- Sort:** alphabetical
- Column list:**
  - amrfinderplus\_amr\_genes
  - amrfinderplus\_amr\_report
  - amrfinderplus\_db\_version
  - amrfinderplus\_stress\_genes
  - amrfinderplus\_stress\_report
  - amrfinderplus\_version
  - amrfinderplus\_virulence\_genes
  - amrfinderplus\_virulence\_report
  - ani\_highest\_percent
  - ani\_highest\_percent\_bases\_aligned
  - ani\_mummer\_version
  - ani\_output\_tsv
  - ani\_top\_species\_match
- Save this column selection** (text)
- Column selection name:** qc\_metrics
- This column selection will be shared with all users of this workspace.
- SAVE** button
- Save this column selection** (callout box)
- CANCEL** and **DONE** buttons at the bottom right.

**Modifying an existing column selection**

The two screenshots show the 'Select columns' dialog box in two states:

- Left Screenshot:** Shows the 'Save this column selection' dialog with the 'Column selection name' field containing 'qc\_metrics'. Below it, a list of 'Your saved column selections' includes 'pulsenet\_genotyping' and 'qc\_metrics'.
- Right Screenshot:** Shows the 'Update' dialog where the 'Column selection name' field is 'qc\_metrics' with a clear button (X). Below it, the 'Your saved column selections' list is the same as in the left screenshot.

Both screenshots include the same column list as the first screenshot, with 'ani\_highest\_percent' and 'ani\_top\_species\_match' selected. The 'UPDATE' button is visible in the right screenshot.

**Appendix PNID01-4a. PulseNet Critical (Pass/Fail) Quality Metrics for Routine Sequence Submissions**

<b>Organism</b>	<b>Average denovo coverage<sup>1</sup></b>	<b>Average quality (Q score)<sup>2</sup></b>	<b>Assembly length (MB)</b>	<b>Secondary species abundance</b>
<i>Listeria monocytogenes</i>	≥ 20x	≥ 30	2.8-3.2	≤ 0.01
<i>E. coli</i> (most serotypes)	≥ 40x	≥ 30	4.9-6.0	≤ 0.01
<i>Shigella spp./Rare E. coli</i>	≥ 40x	≥ 30	4.2-4.9	≤ 0.01
<i>Salmonella spp.</i>	≥ 30x	≥ 30	4.4-5.7	≤ 0.01
<i>Campylobacter spp.</i>	≥ 20x	≥ 30	1.4-2.2	≤ 0.01
<i>Vibrio cholerae</i>	≥ 40x	≥ 30	3.8-4.3	≤ 0.01
<i>Vibrio parahaemolyticus</i>	≥ 40x	≥ 30	4.9-5.5	≤ 0.01
<i>Vibrio vulnificus</i>	≥ 40x	≥ 30	4.7-5.3	≤ 0.01

<sup>1</sup>After quality-based trimming (est\_coverage\_clean)

<sup>2</sup>Before trimming (combined\_mean\_q\_raw)

**Appendix PNID01-4b. TheiaProk Read Pre-Screening Step to Exclude Poor Quality Sequences to Conserve Computational Resources**

The screen task ensures the quantity of sequence data is sufficient to undertake genomic analysis. It uses bash commands for quantification of reads and base pairs, and mash sketching to estimate the genome size and its coverage. At each step, the results are assessed relative to pass/fail criteria and thresholds that may be defined by optional user inputs.

Samples that do not meet these criteria will not be processed further by the workflow:

1. Total number of reads: A sample will fail the read screening task if its total number of reads is less than or equal to min\_reads.
2. The proportion of basepairs reads in the forward and reverse read files: A sample will fail the read screening if fewer than min\_proportion basepairs are in either the reads1 or read2 files.
3. Number of basepairs: A sample will fail the read screening if there are fewer than min\_basepairs basepairs
4. Estimated genome size: A sample will fail the read screening if the estimated genome size is smaller than min\_genome\_size or bigger than max\_genome\_size.
5. Estimated genome coverage: A sample will fail the read screening if the estimated genome coverage is less than the min\_coverage.

Default values:

Int min\_reads = 7472

Int min\_basepairs = 2241820

Int min\_genome\_length = 100000

Int max\_genome\_length = 18040666

Int min\_coverage = 10

Int min\_proportion = 40

**PULSENET INTERNATIONAL STANDARD OPERATING PROCEDURE FOR ANALYZING ILLUMINA SHORT READ WGS DATA USING THE TERRA.BIO PLATFORM**

**Doc. No. PNID01**

**Ver. No. 01**

**Effective Date:**

**Page 50 of 61**

The “raw\_read\_screen” column under “QC\_metrics” will give details if the sample fails the read pre-screen:

The screenshot shows the Terra Bio platform interface. The top navigation bar includes 'WORKSPACES' and 'Data'. Below this, there are tabs for 'DASHBOARD', 'DATA', 'ANALYSES', 'WORKFLOWS', and 'JOB HISTORY'. The 'DATA' tab is active, displaying a table of analysis results. The table has columns for 'analysis\_pt\_24\_id', 'n50\_value', 'number\_contigs', 'raw\_read\_screen', and 'seqs2\_predicted\_contamina...'. The 'raw\_read\_screen' column is circled in red. The table contains several rows of data, including entries for '2011V-1043\_FLEX\_300\_Vibrio', '2012V-1116\_FLEX\_300\_Vibrio', '2013L-5361\_FLEX\_300\_LM', '2013L-5410\_FLEX\_300\_LM', '2013L-5547\_FLEX\_300\_LM', '2013L-5615TK\_NextSeq\_400MB', '2015AM-1304', '2015AM-1305', and '2015C-3794\_FLEX\_300\_Shigella'. The 'raw\_read\_screen' column contains values like 'PASS' and 'FAIL; the estimated coverage is less than the minimum of 10x'.

analysis_pt_24_id	n50_value	number_contigs	raw_read_screen	seqs2_predicted_contamina...
2011V-1043_FLEX_300_Vibrio	124949	77	PASS	
2012V-1116_FLEX_300_Vibrio	498045	39	PASS	
2013L-5361_FLEX_300_LM	526928	12	PASS	
2013L-5410_FLEX_300_LM	526025	15	PASS	
2013L-5547_FLEX_300_LM	435363	20	PASS	
2013L-5615TK_NextSeq_400MB			FAIL; the estimated coverage is less than the minimum of 10x	
2015AM-1304	443204	15	PASS	yes
2015AM-1305	728098	16	PASS	no
2015C-3794_FLEX_300_Shigella	25104	361	PASS	

## Appendix PNID01-5. Customization of a Data Table View for PulseNet Genotyping Assays

**NOTE:** *this customization can only be done after you have run the TheiaProk workflow once.*

1. In the “Data” tab, select the data table of interest, e.g., “CDC\_ATCC\_Sequences”, then select “Settings”.
2. Under “Select columns” the following genotyping assays should be checked:
  - a. Amrfinderplus\_amr\_classes
  - b. Amrfinderplus\_amr\_core\_genes
  - c. Amrfinderplus\_amr\_subclasses
  - d. Amrfinderplus\_virulence\_genes
  - e. Plasmidfinder\_plasmids
  - f. Seqsero2\_predicted\_antigenic\_profile
  - g. Seqsero2\_predicted\_serotype
  - h. Serotypefinder\_serotype
  - i. Ts\_mlst\_predicted\_st

**Select columns**

Show: all | none      Sort: alphabetical

- resfinder\_results
- resfinder\_seqs
- seq\_platform
- seqsero2\_predicted\_antigenic\_profile
- seqsero2\_predicted\_contamination
- seqsero2\_predicted\_serotype
- seqsero2\_report
- seqsero2\_version
- serotypefinder\_docker
- serotypefinder\_report
- serotypefinder\_serotype
- shovill\_pe\_version
- sister\_allele\_fasta
- sistr\_allele\_ison

SAVE THIS COLUMN SELECTION

Your saved column selections:

- pulsenet\_genotyping ⓘ
- qc\_metrics ⓘ

CANCEL      **DONE**

3. Click “Save this column selection”, name the column selection “pulsenet\_genotyping” and click “Save” and then click “Done”.

**NOTE:** *If you are adding or deleting columns from an existing column selection, click “Save this column selection”, select the name from the drop-down menu and click “Update”.*

**PULSENET INTERNATIONAL STANDARD OPERATING PROCEDURE FOR ANALYZING ILLUMINA SHORT READ WGS DATA USING THE TERRA.BIO PLATFORM**

**Doc. No. PNID01**

**Ver. No. 01**

**Effective Date:**

**Page 52 of 61**

**Select columns**

Show: all | none      Sort: alphabetical

resfinder\_pointfinder\_results  
 resfinder\_results  
 resfinder\_seqs  
 seq\_platform  
 seqsero2\_predicted\_antigenic\_profile  
 seqsero2\_predicted\_contamination  
 seqsero2\_predicted\_serotype  
 seqsero2\_report  
 seqsero2\_version  
 serotypefinder\_docker  
 serotypefinder\_report  
 serotypefinder\_serotype  
 shovill\_pe\_version

Save this column selection

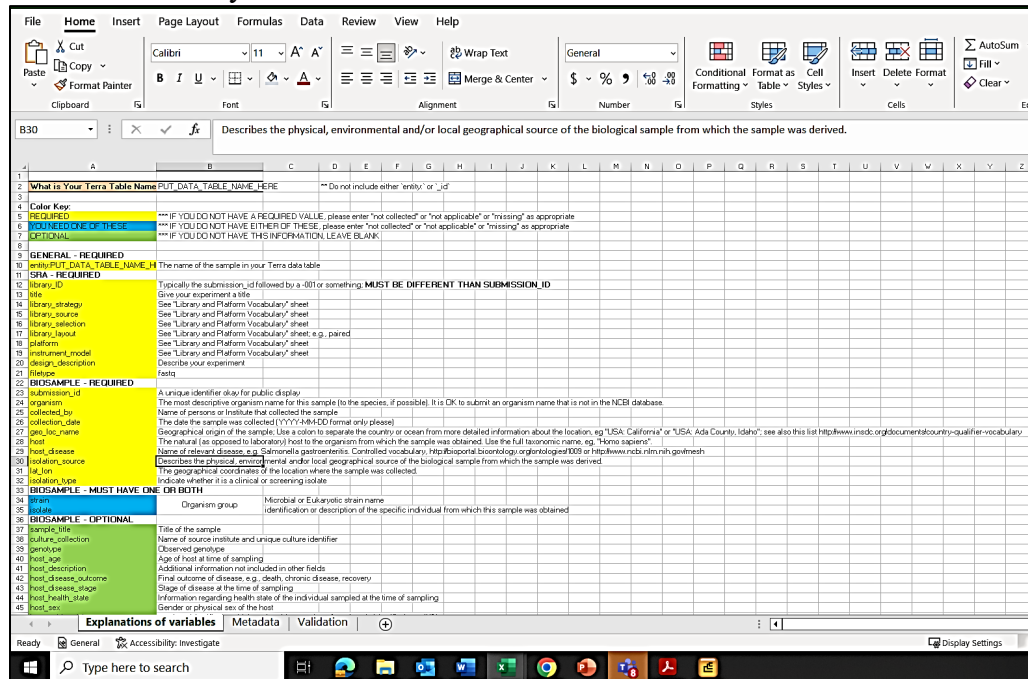
**Column selection name**

Save this column selection      selections:

qc\_metrics

## Appendix PNID01-6. Uploading Additional Metadata to Terra for NCBI Submissions and Customization of a Data Table View for Metadata

- NCBI requires minimal metadata to be uploaded to NCBI to create BioSamples for the sequences to be uploaded to SRA.
- Because the NCBI metadata needs to be formatted in a certain way, please use the **Pathogen** metadata template spreadsheet provided by Theiagen to upload metadata to Terra: [https://theiagen.notion.site/Terra\\_2\\_NCBI-8f014c73acc44465a3d69cf4df93adfe](https://theiagen.notion.site/Terra_2_NCBI-8f014c73acc44465a3d69cf4df93adfe).
- The metadata template has three tabs:
  - The first tab, called “Explanation of variables”, contains descriptions for the required and optional fields.
  - The metadata to be uploaded is entered in the second “Metadata” tab.
  - The third “Validation” tab can be used to ascertain that the required metadata is filled out correctly.



### To upload the Pathogen metadata file to Terra:

1. Fill out the required and optional (if applicable) fields in the metadata tab of the Pathogen metadata template spreadsheet:

**NOTE1:** you can enter “Missing” for any required information you don’t have or you don’t wish to publicly disclose.

**NOTE2:** The metadata below are the PulseNet USA minimum metadata requirements for NCBI uploads designed to protect patient privacy and the integrity of on-going outbreak investigations. On the other hand, enough epidemiologically useful information is provided particularly on non-clinical samples to facilitate attribution.

**PULSENET INTERNATIONAL STANDARD OPERATING PROCEDURE FOR ANALYZING ILLUMINA SHORT READ WGS DATA USING THE TERRA.BIO PLATFORM**

**Doc. No. PNID01**

**Ver. No. 01**

**Effective Date:**

**Page 54 of 61**

- a. In cell A1, enter the data table name into which metadata will be uploaded:  
**entity: data\_table\_name\_id**, e.g., entity: quality\_control\_id.
- b. Entity: enter the sample IDs that match the IDs for the sequences in Terra.
- c. Submission\_id: enter a unique ID for the sample. This is the ID that will be displayed on NCBI.
  - i. Enter the submission\_id also in the “strain” column.
  - ii. Enter the “submission\_id-001” to the “library\_id” column.
- d. Title: “PulseNet”.
- e. For “library\_strategy”, “library\_source”, “library\_selection”, “platform”, “instrument\_model” and “filetype”, pick the correct option from the drop-down menus.
- f. Library layout: “Paired”.
- g. Organism: genus and species.
- h. Collected\_by: the laboratory submitting the sequence.
- i. Collection\_date: **year for clinical isolates, year and month for non-clinical isolates**. Required format: YYYY:MM:DD, e.g., for a clinical sample isolated in 2024: 2024-01-01. For non-clinical sample isolated in February 2024: 2024-02-01.
- j. Geo\_loc\_name: **source country for clinical isolates, source country and state (or other more detailed location) for non-clinical isolates**. Use a colon to separate the country and a more detailed location, e.g., USA:CA.
- k. Isolation\_source: **“Missing” for clinical isolates, the exact source for non-clinical isolates**, e.g., lettuce, chicken breast, swab, etc.
- l. Isolation\_type: clinical, environmental, food or animal.
- m. Serotype: *E. coli* serotype.
- n. Serovar: *Salmonella* serovar.

***Clinical (human origin) sequences only***

	A	B	C	D	E	F	G	H	I	J	K	L
1	entity:quality_control_id	library_ID	title	library_strategy	library_source	library_selection	library_layout	platform	instrument_model	design_description	filetype	submission_id
2	2017C-4936	2017C-4936-001	PulseNet	WGS	GENOMIC	RANDOM	paired	ILLUMINA	Illumina HiSeq 2500	Missing	fastq	2017C-4936
3	2018C-4039	2018C-4039-001	PulseNet	WGS	GENOMIC	RANDOM	paired	ILLUMINA	Illumina HiSeq 2500	Missing	fastq	2018C-4039
4	2019C-3204	2019C-3204-001	PulseNet	WGS	GENOMIC	RANDOM	paired	ILLUMINA	Illumina HiSeq 2500	Missing	fastq	2019C-3204

	K	L	M	N	O	P	Q	R	S	T	U	V	W	X
1	letype	submission_id	organism	collected_by	collection_date	geo_loc_name	host	host_disease	isolation_source	lat_lon	isolation_type	strain	isolate	sample_title
2	fastq	2017C-4936	Escherichia coli	CDC	2017-01-01	USA	Homo sapiens	Missing	Missing	Missing	clinical	2017C-4936		
3	fastq	2018C-4039	Escherichia coli	CDC	2018-01-01	USA	Homo sapiens	Missing	Missing	Missing	clinical	2018C-4039		
4	fastq	2019C-3204	Shigella sonnei	CDC	2018-01-01	USA	Homo sapiens	Missing	Missing	Missing	clinical	2019C-3204		

**PULSENET INTERNATIONAL STANDARD OPERATING PROCEDURE FOR ANALYZING ILLUMINA SHORT READ WGS DATA USING THE TERRA.BIO PLATFORM**

**Doc. No. PNID01**

**Ver. No. 01**

**Effective Date:**

**Page 55 of 61**

*Clinical (human origin) and non-clinical sequences*

2. Switch to the “Validation” tab to confirm that the metadata template has been filled out correctly.

**NOTE:** *the validation of the proper format of the geological location will only pass if the source country and a more detailed location are both specified. Therefore, validation will*

**PULSENET INTERNATIONAL STANDARD OPERATING PROCEDURE FOR ANALYZING ILLUMINA SHORT READ WGS DATA USING THE TERRA.BIO PLATFORM**

**Doc. No. PNID01**

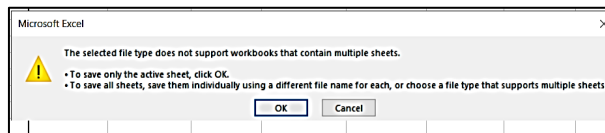
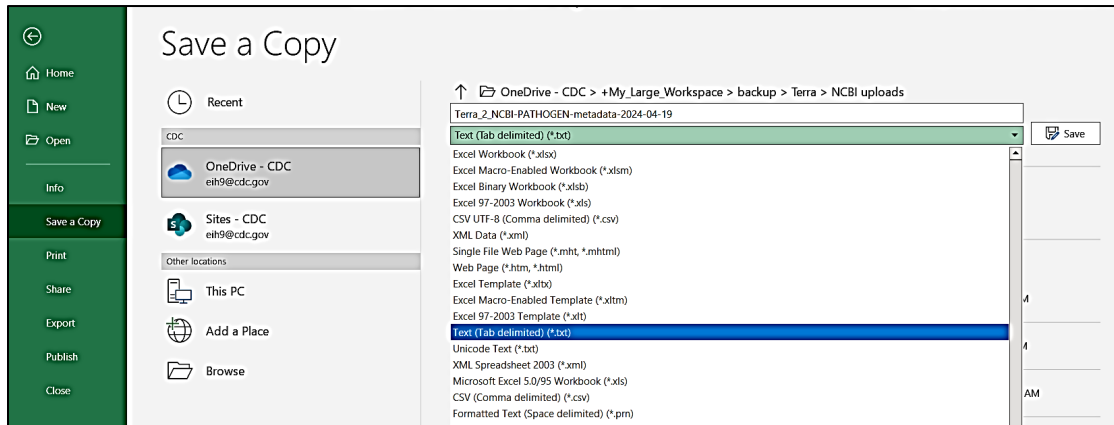
**Ver. No. 01**

**Effective Date:**

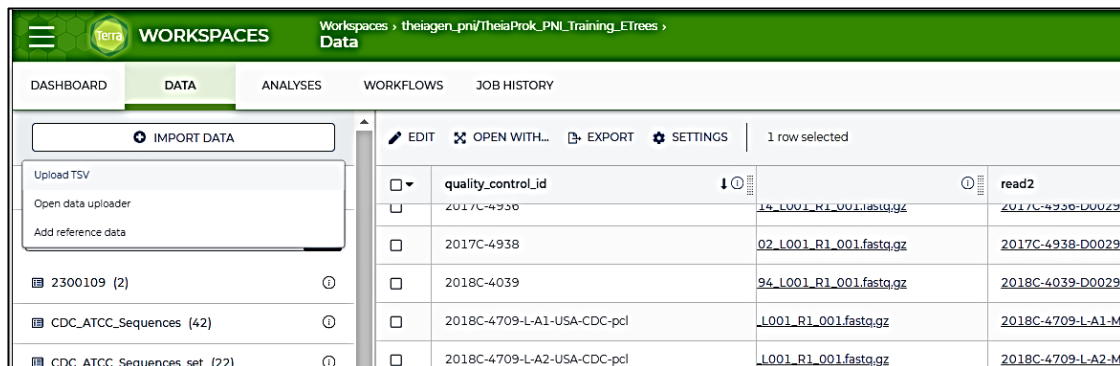
**Page 56 of 61**

*not pass for clinical isolates for which only the source country is specified. However, NCBI will accept source country as the only geological location.*

3. Save the filled-out metadata template as a tsv (tab delimited) file. Click “OK” on the pop-up window stating that the selected file type does not support workbooks with multiple sheets.



4. In the Terra Workspaces “Data” tab, click on “Import Data” and select “Upload TSV” from the drop-down menu.



5. In the “Import Table Data” pop-up window “File import” tab, click in the middle to select the tsv file.

**Import Table Data**

Choose the data import option below. [Click here for more info on the table.](#)

Data will be saved in Terra-managed location: us-us-central1 (Iowa)

**FILE IMPORT**    TEXT IMPORT

---

Select the **TSV** file containing your data:

Drag or Click to select a .tsv file

Selected File: **None**

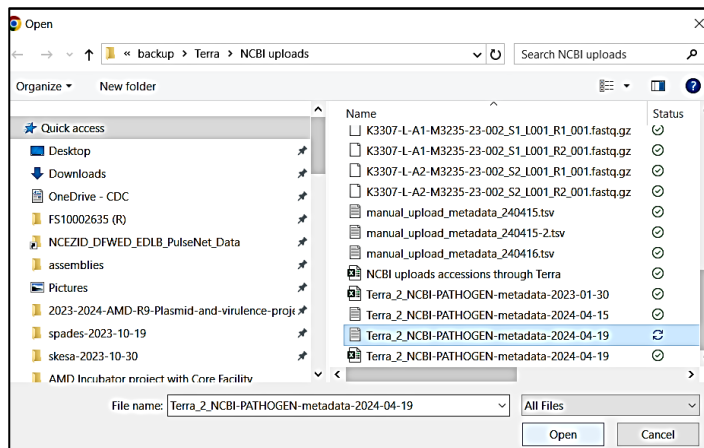
---

**TSV file templates**

[Download sample\\_template.tsv](#)

Terra Support: Importing Data - Using a Template

- Navigate to the location where the metadata tsv file is saved, select the file and click “Open”.



- In the “Import Table Data” pop-up window, you will get a warning that data already exists in the data table in question and uploading more data to it may overwrite existing data. Also, there will be a warning if the metadata tsv file does not have the same information (some data missing for some sequences) for all entries. Click on “Start Import Job”.

**Import Table Data**

Choose the data import option below. [Click here for more info on the table.](#)

Data will be saved in Terra-managed location: us-us-central1 (Iowa)

**FILE IMPORT**    TEXT IMPORT

---

Select the **TSV** file containing your data: Clear

Drag or Click to select a .tsv file

Selected File: Terra\_2\_NCBI-PATHOGEN-metadata-2024-04-19.txt

**Warning:** Data with the type 'quality\_control' already exists in this workspace. Uploading more data for the same type may overwrite some entries.

**Warning:** We have detected empty cells in your TSV. Please choose an option:

Ignore empty cells (default)

Overwrite existing cells with empty cells

---

**TSV file templates**

[Download sample\\_template.tsv](#)

Terra Support: Importing Data - Using a Template

**PULSENET INTERNATIONAL STANDARD OPERATING PROCEDURE FOR ANALYZING ILLUMINA SHORT READ WGS DATA USING THE TERRA.BIO PLATFORM**

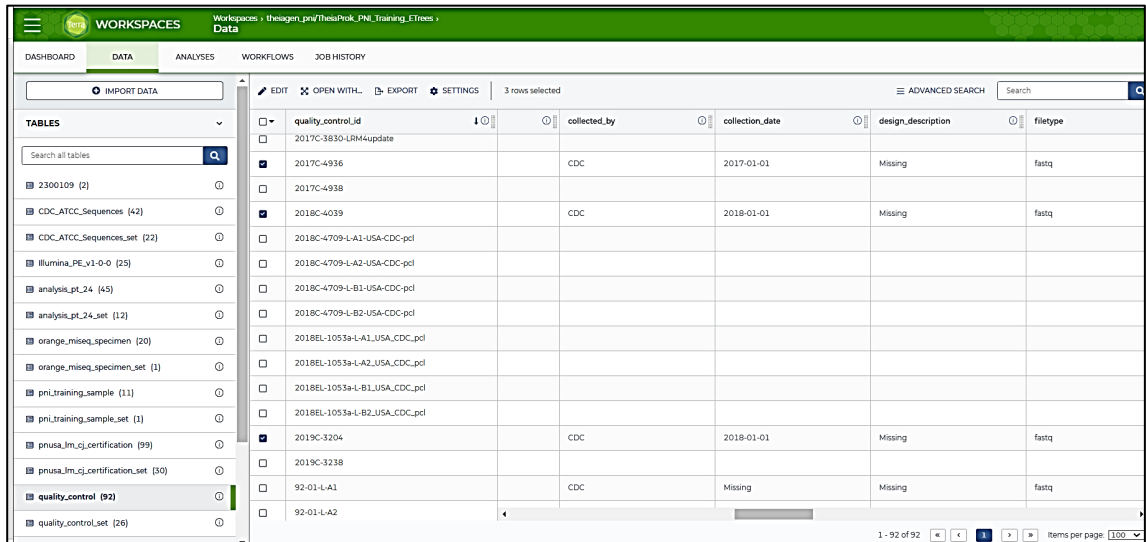
**Doc. No. PNID01**

**Ver. No. 01**

**Effective Date:**

**Page 58 of 61**

8. After the upload is done you should see the desired metadata fields populated in the data table.

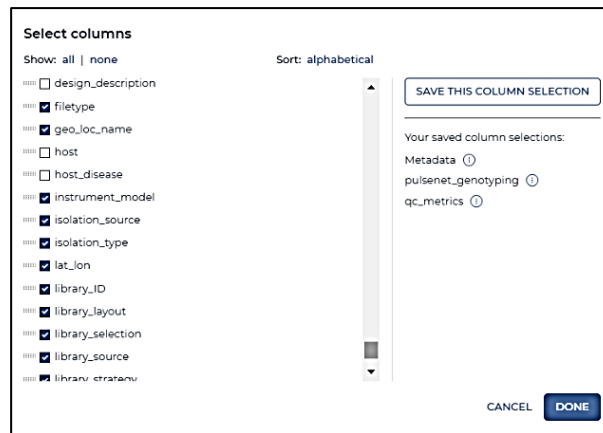


quality_control_id	collected_by	collection_date	design_description	filetype
2017C-3830-LRH4update				
2017C-4936	CDC	2017-01-01	Missing	fastq
2017C-4938				
2018C-4039	CDC	2018-01-01	Missing	fastq
2018C-4709-L-A1-USA-CDC-pcd				
2018C-4709-L-A2-USA-CDC-pcd				
2018C-4709-L-B1-USA-CDC-pcd				
2018C-4709-L-B2-USA-CDC-pcd				
2018EL-1053a-L-A1-USA-CDC-pcd				
2018EL-1053a-L-A2-USA-CDC-pcd				
2018EL-1053a-L-B1-USA-CDC-pcd				
2018EL-1053a-L-B2-USA-CDC-pcd				
2019C-3204	CDC	2018-01-01	Missing	fastq
2019C-3238				
92-01-L-A1	CDC	Missing	Missing	fastq
92-01-L-A2				

**Create a separate metadata view for the data table**

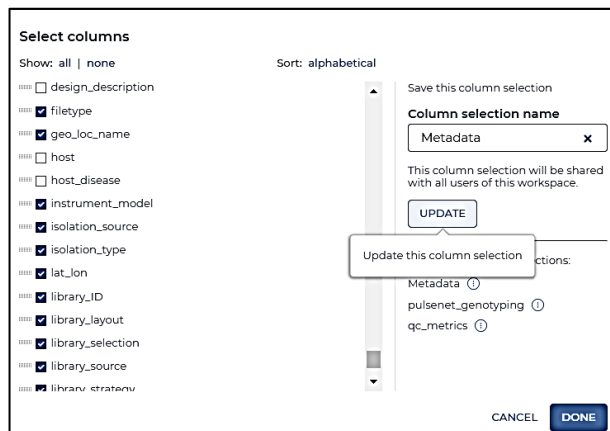
1. In the “Data” tab, select the data table of interest, then select “Settings”.
2. Under “Select columns”, check all the desired metadata columns.
  - a. Recommended metadata for NCBI submissions:
    - i. collected\_by
    - ii. collection\_date
    - iii. filetype
    - iv. geo\_loc\_name
    - v. instrument\_model
    - vi. isolation\_source
    - vii. isolation\_type
    - viii. library\_id
    - ix. library\_layout
    - x. library\_selection
    - xi. library\_source
    - xii. library\_strategy
    - xiii. organism
    - xiv. platform
    - xv. strain
    - xvi. submission\_id
    - xvii. title
    - xviii. serotype
    - xix. serovar
  - b. Additional useful information about the sequence
    - i. read1 (R1 FASTQ file name)
    - ii. read2 (R2 FASTQ file name)
    - iii. assembly\_fasta (location for the assembly generated by Terra)

- iv. If uploading data directly from the Illumina BaseSpace:
  - 1. basespace\_collection\_id
  - 2. basespace\_fetch\_analysis\_date
  - 3. basespace\_fetch\_version
  - 4. basespace\_sample\_id
  - 5. basespace\_sample\_name
- v. biosample\_accession
- vi. sra\_accession



- 3. Click on “Save this column selection”.
- 4. Name the column selection “Metadata” and click on “Save” and “Done”.

**NOTE:** *If you are adding or deleting columns from an existing column selection, click “Save this column selection”, select the name from the drop-down menu and click “Update”.*



- 5. Only the desired metadata columns should now be visible in the data table

# PULSENET INTERNATIONAL STANDARD OPERATING PROCEDURE FOR ANALYZING ILLUMINA SHORT READ WGS DATA USING THE TERRA.BIO PLATFORM

Doc. No. PNID01

Ver. No. 01

Effective Date:

Page 60 of 61

The screenshot shows the Terra Bio Workspaces interface. The top navigation bar includes 'WORKSPACES' and 'Data'. Below this, there are tabs for 'DASHBOARD', 'DATA', 'ANALYSES', 'WORKFLOWS', and 'JOB HISTORY'. The 'DATA' tab is active, displaying a table with columns: 'quality\_control\_id', 'pile\_accession', 'collected\_by', 'collection\_date', 'filetype', and 'geo.'. The table contains 14 rows of data, including quality control records for various projects like 'orange\_miseq...', 'pni\_training\_sa...', 'pnusa\_fm\_cj\_ce...', and '2019C-3204'. A left sidebar shows a list of workspace folders such as 'analysis\_pt\_24', 'orange\_miseq...', 'pni\_training\_sa...', 'pnusa\_fm\_cj\_ce...', 'quality\_control', 'REFERENCE DATA', and 'OTHER DATA'. The bottom of the interface shows pagination information: '1 - 94 of 94' and 'Items per page: 100'.

quality_control_id	pile_accession	collected_by	collection_date	filetype	geo.
2017C-4936	1039458	CDC	2017-01-01	fastq	USA
2017C-4938					
2018C-4039	1039457	CDC	2018-01-01	fastq	USA
2018C-4709-L-A1-USA-CDC-pcl					
2018C-4709-L-A2-USA-CDC-pcl					
2018C-4709-L-B1-USA-CDC-pcl					
2018C-4709-L-B2-USA-CDC-pcl					
2018EL-1053a-L-A1_USA_CDC_pcl					
2018EL-1053a-L-A2_USA_CDC_pcl					
2018EL-1053a-L-B1_USA_CDC_pcl					
2018EL-1053a-L-B2_USA_CDC_pcl					
2019C-3204	1039456	CDC	2017-01-01	fastq	USA
2019C-3238					

**Appendix PNID01-7: Overview of the TheiaProk Workflow for Bacterial Characterization**

