

2 ▾

Oct 18, 2021

Auger/Auspice/UShER SARS-CoV-2 Cluster Detection Workflow for the Terra Platform V.2

protocol



1

Anusha Ginni¹, Jill V Hagey¹, Michael Weigand¹,
Technical Outreach and Assistance for States Team¹

¹Centers for Disease Control and Prevention

1



protocol .

Anusha Ginni

The opinions expressed here do not necessarily reflect the opinions of the Centers for Disease Control and Prevention or the institutions with which the authors are affiliated. The protocol content here is under development and is for informational purposes only and does not constitute legal, medical, clinical, or safety advice, or otherwise; content added to [protocols.io](#) is not peer reviewed and may not have undergone a formal approval of any kind. Information presented in this protocol should not substitute for independent professional judgment, advice, diagnosis, or treatment. Any action you take or refrain from taking using or relying upon the information presented here is strictly at your own risk. You agree that neither the Company nor any of the authors, contributors, administrators, or anyone else associated with [protocols.io](#), can be held responsible for your use of the information contained in or linked to this protocol or any of our Sites/Apps and Services.

The Titan_Augur workflow processes concatenated assembly fasta data for SARS-CoV-2 phylogenetic analysis and cluster detection using Augur, the modular bioinformatics engine behind Nextstrain. Two workflows, Titan_Augur_Prep and Titan_Augur_Run, must be completed in series and produce the required inputs for visualization with Auspice to evaluate genetic relatedness among sequences or draw phylogenetic inference.

Inputs for Titan_Augur_Prep include the assembled fasta, metadata (uploaded tsv), and PANGO lineage assignments for each sample. Running Titan on Terra generates the necessary assembly fasta and PANGO lineage assignments, but these data can also be uploaded if prepared using another pipeline. Titan_Augur_Prep processes these data and outputs the correctly formatted metadata table and concatenated fasta files required for subsequent Titan_Augur_Run analyses. Titan_Augur_Run then generates the inputs needed for phylogenetic inference on the Auspice or Nextstrain interactive visualization tools.

This protocol also includes steps for visualization using Auspice as well as the UCSC UShER webportal.

For technical assistance, please contact: TOAST@cdc.gov

Anusha Ginni, Jill V Hagey, Michael Weigand, Technical Outreach and Assistance for States Team 2021. Auger/Auspice/UShER SARS-CoV-2 Cluster Detection Workflow for the Terra Platform. [protocols.io](#)
<https://protocols.io/view/auger-auspice-usher-sars-cov-2-cluster-detection-w-bymepu3e>
Technical Outreach and Assistance for States Team

(i)



(ILLUMINA MENU) Protocols for SARS-CoV-2 Library Prep with ARTIC Primers, Bioinformatic Analysis, and Database Submission

Nanopore, SARS-CoV-2, Pangolin, Genomics, Virology, RNA, Covid, Computational Biology, Sequencing, Phylogenetics, Auspice, NextStrain, GISAID, UShER, Augur

protocol ,

Sep 28, 2021

Oct 18, 2021

Oct 18, 2021 | Anusha Ginni

53638

Part of collection

(ILLUMINA MENU) Protocols for SARS-CoV-2 Library Prep with ARTIC Primers, Bioinformatic Analysis, and Database Submission

:

The opinions expressed here do not necessarily reflect the opinions of the Centers for Disease Control and Prevention or the institutions with which the authors are affiliated. The protocol content here is under development and is for informational purposes only and does not constitute legal, medical, clinical, or safety advice, or otherwise; content added to [protocols.io](#) is not peer reviewed and may not have undergone a formal approval of any kind. Information presented in this protocol should not substitute for independent professional judgment, advice, diagnosis, or treatment. Any action you take or refrain from taking using or relying upon the information presented here is strictly at your own risk. You agree that neither the Company nor any of the authors, contributors, administrators, or anyone else associated with [protocols.io](#), can be held responsible for your use of the information contained in or linked to this protocol or any of our Sites/Apps and Services.

Setup Terra and Google Cloud Accounts

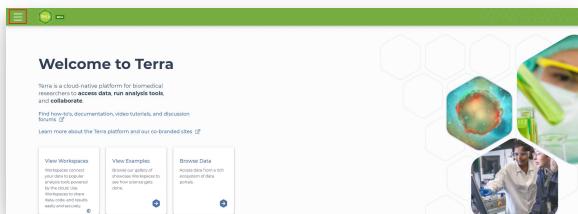
1

protocols.io

2

This is an open access protocol distributed under the terms of the **Creative Commons Attribution License** (<https://creativecommons.org/licenses/by/4.0/>), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited

The Terra platform registration requires a Google account. If you have a Google account you can sign in using the Terra login page:
<https://app.terra.bio/>



Welcome page for Terra.bio.

If you do not have Google email you can set up a Google account with a non-Google email. The steps to do this are described in the following link:

<https://support.terra.bio/hc/en-us/articles/360029186611>

The Terra platform uses the Google Cloud to run workflows and store data. The following documentation will describe how to set up a Google Cloud account:

<https://support.terra.bio/hc/en-us/articles/360046295092>

To link your Terra platform account with your Google Cloud account follow the instructions provided in the following link:

<https://support.terra.bio/hc/en-us/articles/360026182251-How-to-set-up-billing-in-Terra>

Scroll down to Section 3 titled "Create a Terra Billing Project" and follow the instructions. **It is important to note that the name you give to your Terra Billing Project must be unique across all Google Billing Projects.** If the name provided is not unique, it won't immediately throw an error, but instead will not complete the process of associating the Google Billing account with the Terra Billing Project. If this occurs cancel the process and set up a new Terra Billing Project.

For detailed step by step instructions, refer to Terra registration check out our [protocol on registration and billing](#).

Import Augur workflows from Dockstore

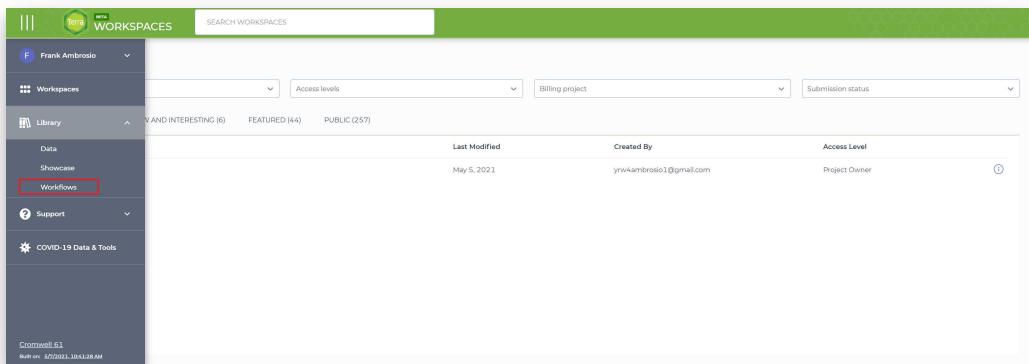
2

Importing the Titan_Augur Workflows from Dockstore to the Terra User Workspace

Note: The Augur workflow includes two steps: *Titan_Augur_Prep* and *Titan_Augur_Run*.

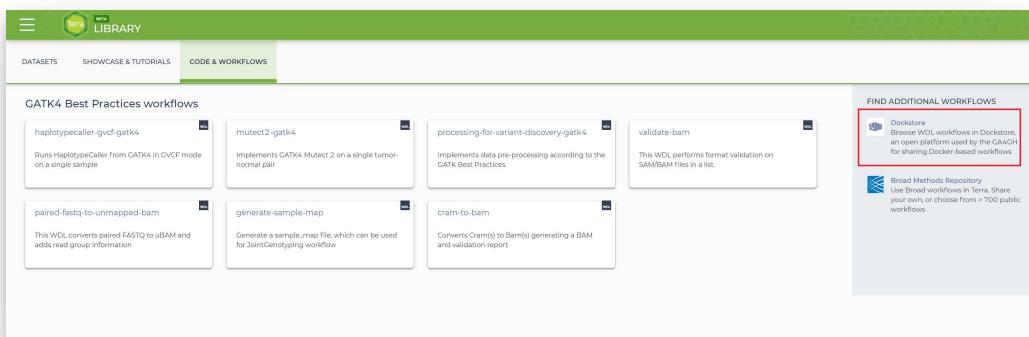
The *Titan_Augur_Prep* workflow prepares the metadata and assembly fasta required by the *Titan_Augur_Run* workflow, so both pipelines must be imported into your workspace.

The Augur workflows are hosted on the [Dockstore repository](#) and must be imported into the user's Terra Workspace. Begin by clicking on the three parallel lines in the top left-hand corner, followed by clicking the 'Library' tab and finally click the 'Workflows' button.



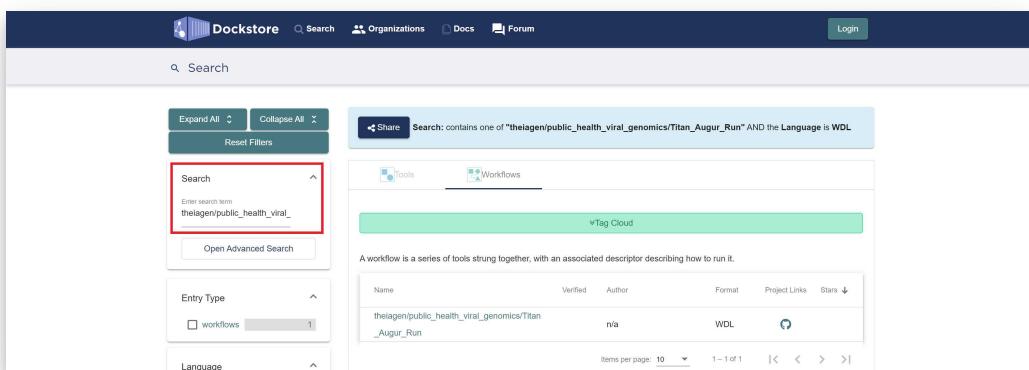
Find the 'Workflows' button listed under the 'Library' tab in the selection panel.

On the right side of the page under 'Find Additional Workflows' select 'Dockstore' from the dialogue box .



Workflows panel with link to Dockstore.

On the left side of the Dockstore page search for 'theiagen/public_health_viral_genomics/Titan_Augur_Run' in the search bar.



Search results for theiagen/public_health_viral_genomics/Titan_Augur_Run

Click the 'theiagen/public_health_viral_genomics/Titan_Augur_Run' link. This will take you to a page where you can import the workflow into your Terra workspace.

Dockstore link to terra

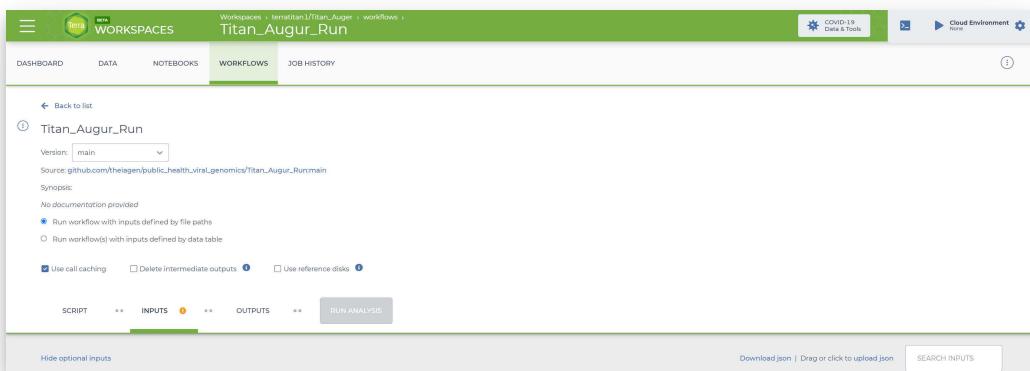
On the right-hand side of the page under the 'Launch with' window click the 'Terra' button. It should bring you back to the Terra platform within the 'Import Workflow' page. If you already have a workspace select it from the dropdown menu Under 'Destination Workspace' and click 'Import'. If not you will need to create a workspace first (see note below)

If you need to create a new workspace click the 'create a new workspace' button.

A pop-up window titled 'Create a New Workspace' should appear. Name your new workspace and associate it with a billing account using the 'Billing project' drop-down menu (this is the Terra Billing Account you created in the previous step should be available). Finally, click the 'Create Workspace' button.

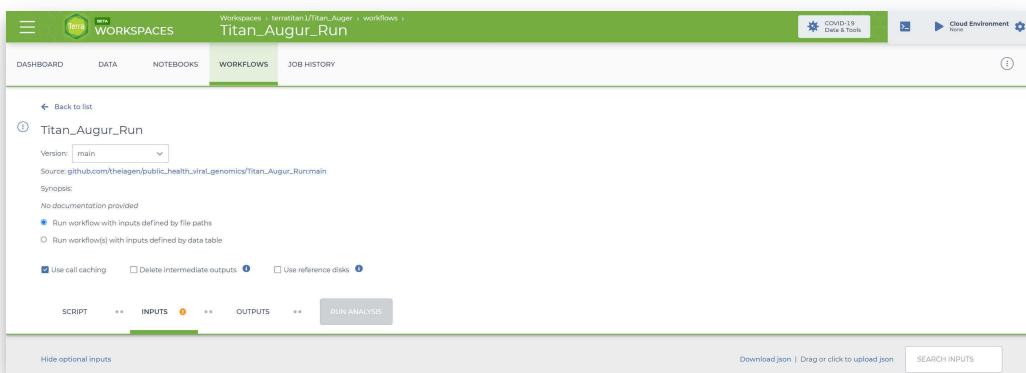
New Workspace Panel

After clicking the 'Create Workspace' button you should be automatically directed to the Augur workflow panel in the new workspace page that was just created.



If you have imported the workflow, you should see this page.

After clicking the 'Import' button you should be automatically directed to the Augur workflow panel in your specified workspace.



If you have imported the workflow, you should see this page.

Now complete the same steps to import the **theiagen/public_health_viral_genomics/Titan_Augur_Prep** workflow as well.

The video below demonstrates how to import both the Titan_Augur_Prep and Titan_Augur_Run workflows.

Analysis Inputs/Data uploading

3

Before you begin: Datasets with <15 samples must include some genetic diversity and varied

collection dates to ensure the Augur workflow runs successfully. This can be achieved by adding an outlier sample with a discrepant lineage and collection date.

The Augur workflows require the assembled fasta sequence (consensus sequence), PANGO lineage assignments, and metadata for each sample. These required inputs are listed on the workflow page:

Task name	variable	Type	Attribute
titan_augur_prep	assembly	String	Required
titan_augur_prep	collection_date	String	Required
titan_augur_prep	iso_continent	String	Required
titan_augur_prep	iso_country	String	Required
titan_augur_prep	iso_state	String	Required
titan_augur_prep	pango_lineage	String	Required

The Titan_Augur_Prep workflow page where required inputs can be added from the existing workspace.

The assembly file is generated from your genomic characterization workflow. If you have run Titan workflows on Terra then you will already have generated the assembly fasta files and all you need at this point is to select the samples from your analysis sample sets within your workspace. We will do this in step 4, but will still need to upload your metadata for each sample if that is not already associated with them.

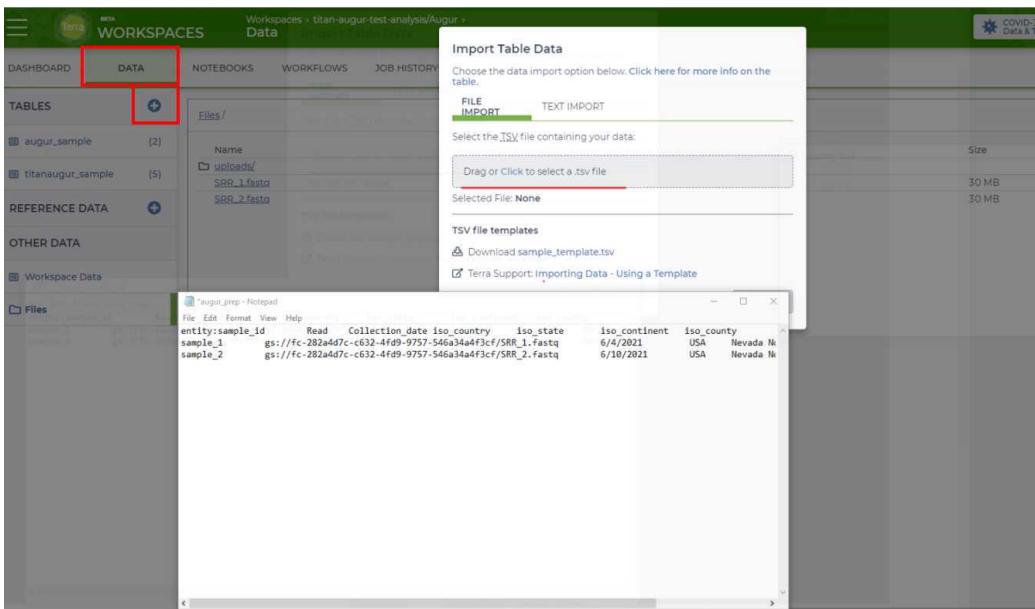
Construct a file with the following columns that are [required metadata](#) for each sample of interest:

- Collection_date
- iso_country
- iso_state
- iso_continent
- iso_county

Here is an example for TABLE file that you will upload to your workspace:

[**augurprep_metadata.tsv**](#)

To associate your metadata with the assembly fasta and pango lineages already in Terra, being by going to the DATA page and click the blue "+" sign next to the TABLES tab on the left of the web page. A pop-up window will appear to import a table file (see example file above) and either drag and drop the file or browse to find it. Once uploaded return to your samples and notice you will now how additional columns associated with each sample. **Terra will know to associate the columns with the correct row by your "entity:" column so use the same sample_ids as you did when you originally uploaded the data.**



Uploading metadata to associate with samples already run through Titan on Terra.

If you have already performed genome characterization on Terra and all the metadata is associated with your samples, skip to Step 4 to select the Titan_Augur_Prep and Titan_Augur_Run required inputs from the TABLE.

If you did not perform any previous genomic characterization of SARS-CoV-2 sequences on Terra, the assemblies and associated metadata file must first be uploaded to run the Augur workflows (see sub-steps below).

You have 3 options for uploading assemblies and metadata. These are the same as what you saw in the Titan protocols.

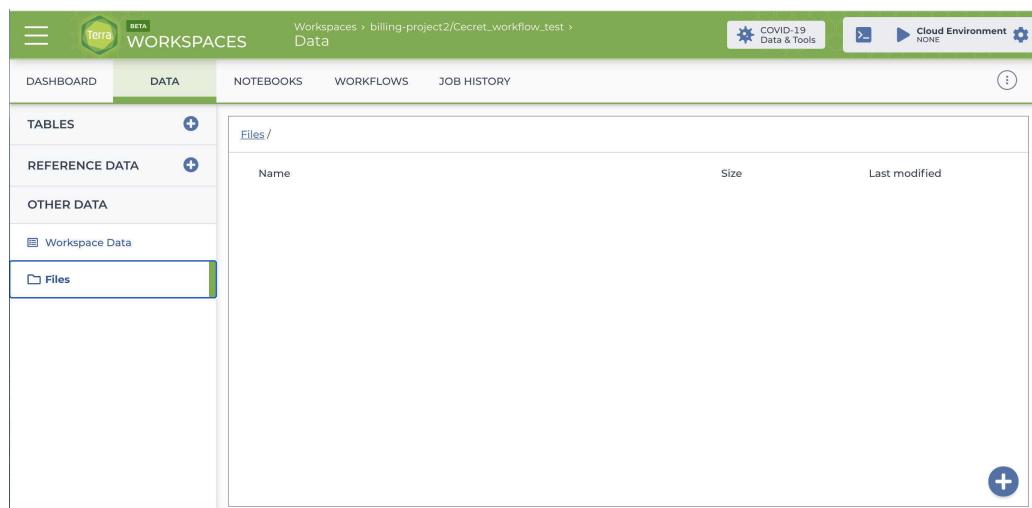
1. Upload on Terra (can only do single files).
2. Upload to Google bucket and link to Terra (can do single files or bulk files/folders).
3. Upload via '<https://app.terra.bio/#upload>' -- **This is the easiest option**

Here we will show you how to select/upload metadata information to associate with your samples.

Note: All the variable names and metadata information in the instructions are dummy and are just for demo purpose only. Please make sure you are using the right variables/names that pertain to your samples.

3.1 The first option is to upload a single sample at a time:

Inside your new workspace page, click on the 'Data' panel in the newly created workspace and then click on the 'Files' tab.



The 'Files' tab within the 'Data' panel of the newly created workspace

Once in the 'Files' tab you can either just drag and drop your files into this space or move the mouse over the blue plus sign icon in the bottom right-hand corner and click 'upload'. Upload the sequence files you'll need for this analysis. **Each fastq sequence file has to be uploaded individually using this method.**

Once the files are uploaded you will need to bring a table in to associate your files with their corresponding link to their google bucket location.

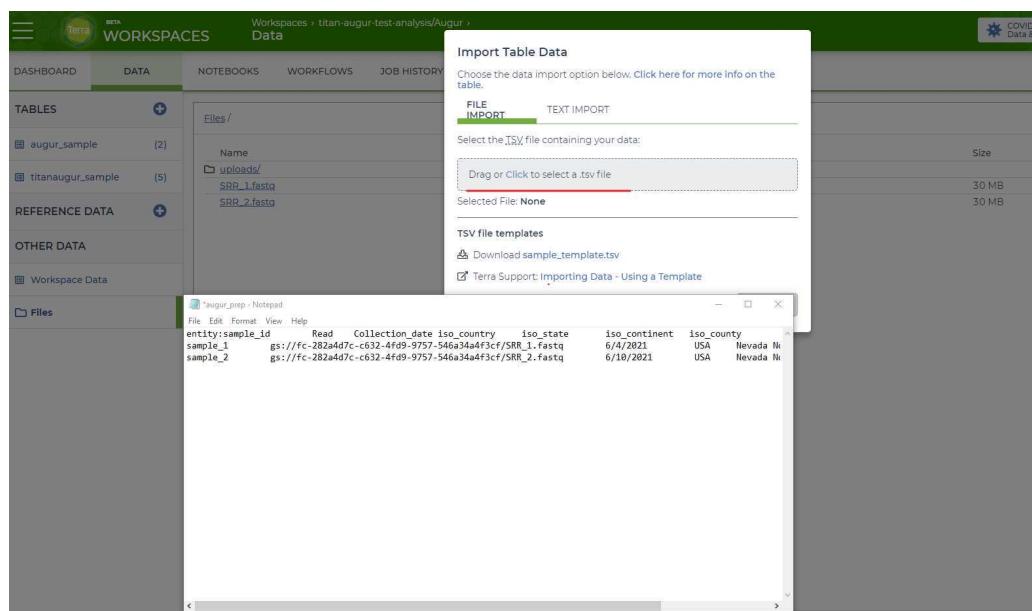


TABLE data updated with the file paths before uploading to workspace.

The Terra sample table file has to follow a specific template. We've provided the template file here [@ Terra_augur_Table_Upload.txt](#) as a downloadable **tab-separated** (or .tsv) file. This example file has the **minimum amount** of columns to be able to create a collection using this method.

The tab-separated table has two columns since we have ONT data in this example:entity:sample_id and fastas. If you had Illumina PE data you would need three columns entity:sample_id, Forward_Read, and Reverse_Read.

Either by editing the text file or using spreadsheet software like Excel, fill in each column with the required information. The first column 'entity:Test_augur_sample_id' is the sample name that is provided by the user. The second column, fasta are the file paths where it is stored within the Google Cloud.

While you had to upload your samples individually you can have all your samples in one datatable.

To identify the Google Cloud location right-click on each fastq file that was uploaded in the previous step and copy the link address. The file path should look like something similar to the following:
gs://fc-b1e3191a-3d9f-43fe-9743-255551ce2f38/SRR11953697.fastq.gz

Once the table is filled in with the required information be sure to save it as a **tab-separated** (or tsv) file. The spreadsheet software should have an option to save as a 'tsv' file.

Important note on column names: DO NOT USE SPACES! As we did before in creating our workspace and collection names use "_" instead of spaces! The first column MUST have the name "entity:sample_id", you can add other things between "entity:" and "sample_id" if you want. Example: "entity:Illumina_sample_id". You can call the other columns whatever you like, but obviously clarity is key.

When completed the table should look similar to the following:

	A	B
1	entity:Test_augur_sample_id	fastas
2	ERR6000261	gs://fc-7312f3f9-5686-4ec3-bcee-51dd2fb5dea2/uploads/augur_data/ERR6000261.fasta
3	ERR6000262	gs://fc-7312f3f9-5686-4ec3-bcee-51dd2fb5dea2/uploads/augur_data/ERR6000262.fasta
4		
5		
6		
7		
8		

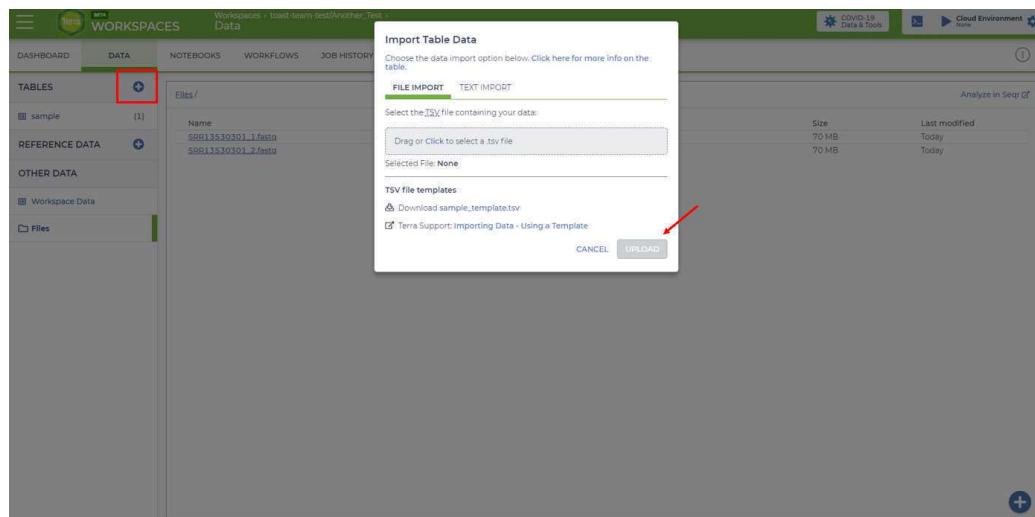
Example Terra sample file

Once the table is filled in with the required information be sure to save it as a **tab-separated** (or tsv) file. The spreadsheet software should have an option to save as a 'tsv' file.

Note: When uploading additional sample table files to the same workspace the entity types must be unique and end in "_id" (e.g. sample1_id, sample2_id etc.)

Finally, the completed Terra sample file will need to be uploaded to the workspace. On

the right-hand side of the workspace 'Data' panel there is a 'Tables' tab. Click the blue plus sign icon on the right edge of the 'Tables' tab. A popup window should appear titled 'Import Table Data'. Select your completed tab-separated sample file for upload and then click the 'Upload' button.



The 'Import Table Data' window for uploading the Terra sample file

If the upload is successful then the sample file should be located under the 'Tables' tab as 'sample (#)' where # is the number of samples in your file.

titanaugur.sample (5)					
	Collection_date	iso_continent	iso_country	iso_county	iso_state
sample_1	6/4/2021	North America	USA	Carson City	Nevada
sample_2	6/10/2021	North America	USA	Clark County	Nevada

The workspace 'Data' panel after successfully uploading the fastq sequence files and TABLE file

The whole process is shown in the video below:

You will notice in the video that you need to upload a file with sample metadata and the path to the fasta/fastq files. Once you have uploaded the files into the workspace, you can either copy, paste and save the paths to the metadata/Table file and upload or edit the paths once you have uploaded the metadata (this way is shown in the video)

- 3.2** You will likely have many samples to upload and you can do this by going directly to your Google bucket.

First, go to “DASHBOARD” tab in your workspace and click “Google Bucket” at the bottom right corner of the same page.

The screenshot shows the NextBIO Workspaces interface. At the top, there's a green header bar with the NextBIO logo, the word "WORKSPACES", and a dropdown menu. Below the header, there are tabs for DASHBOARD, DATA, NOTEBOOKS, WORKFLOWS, and JOB HISTORY. The DASHBOARD tab is selected. On the left, there's a sidebar with "ABOUT THE WORKSPACE" and a message "Just trying it out.". To the right, there's a "WORKSPACE INFORMATION" section with details like Creation Date (3/16/2021), Last Updated (5/7/2021), Subscribers (1), Access Level (Proj. Owner), and Project ID (GOOGLEPROJECTID). Below that is an "OWNERS" section with an email address (apik@cdc.gov) and a "TAGS" section with a placeholder "Add a tag". At the bottom right, there's a "Google Bucket" section with a name (fc-c3f6b8d4-bbed-44a3-9fe6-6e87813941a9) and a link to open it in a browser.

The Dashboard tab of your workspace.

This will direct you to your “Google Cloud Platform” page for data uploading.

The screenshot shows the Google Cloud Platform Cloud Storage interface. The left sidebar has "Cloud Storage" selected. The main area shows "Bucket details" for "fc-c3f6b8d4-bbed-44a3-9fe6-6e87813941a9". There are tabs for OBJECTS, CONFIGURATION, PERMISSIONS, RETENTION, and LIFECYCLE. The OBJECTS tab is selected. It shows a table with columns: Name, Size, Type, Created time, Storage class, Last modified, Public access, Encryption, Retention expiration date, and Hold. There are three entries: a folder named "60e211ce-1f02-4515-80ca-dbf0c440469c/" and a folder named "uploads/". Below the table, there are buttons for "UPLOAD FILES", "UPLOAD FOLDER", "CREATE FOLDER", "MANAGE HOLDS", "DOWNLOAD", and "DELETE".

Click “UPLOAD FILES” in the middle of this page to upload single or multiple fastq files. Or you can click “UPLOAD FOLDER” to upload a folder with multiple fastq files stored inside.

The screenshot shows a file upload interface. On the left, there's a sidebar with 'Quick access' (Desktop, Downloads, fastfiles, Pipeline_Review, Terra, fastq), 'OneDrive - CDC' (This PC, Local Disk (C:)), and 'Network'. A central window displays a file selection dialog with the path 'File name: "SRR_2.fastq" * SRR_1.fastq'. Below it is a large circular icon with a plus sign. To the right, a preview pane shows a list of files under 'fc-282a4d7c-c632-4fd9-9757-546a34a4f3cf' bucket, including 'SRR_1.fastq', 'SRR_2.fastq', 'SRR_3.fastq', 'SRR_4.fastq', and 'SRR_5.fastq'. At the bottom, a message says 'Your bucket is ready. Just add data.' and 'Drop files and folders here or use the upload buttons above. Looking for non-current versions of object? Use [curl](#) or the APIs.'

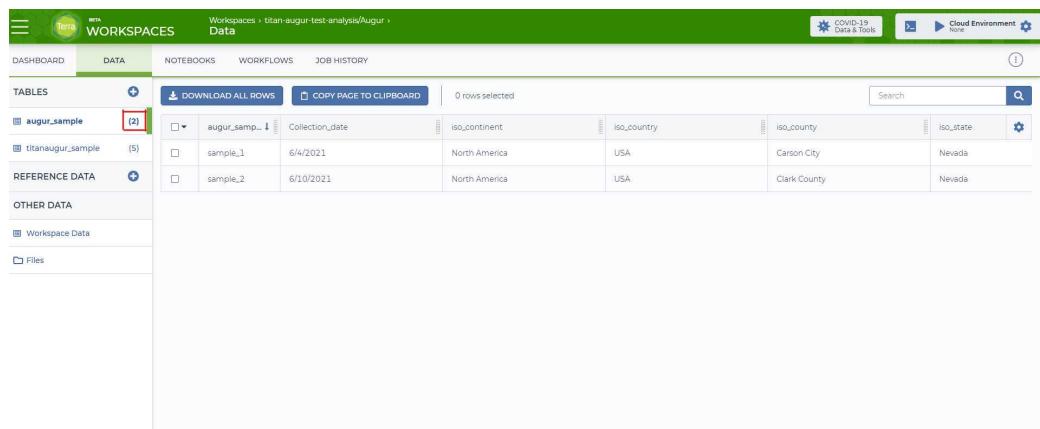
Go back to your Terra account and click “DATA” tab. The successfully uploaded fastq files will show up. If the upload is successful then the files should be located under the ‘Files’ tab, either inside the folder (as you named) or files

The screenshot shows the Terra Data section. The left sidebar includes 'WORKSPACES' (selected), 'DASHBOARD', 'DATA' (selected), 'NOTEBOOKS', 'WORKFLOWS', and 'JOB HISTORY'. Under 'FILES', there are sections for 'TITAN-AUGUR' (with 'augur_sample' and 'titan-augur_sample' subfolders) and 'REFERENCE DATA' (with 'augur' subfolder). The main area shows a table of uploaded files:

Name	Size	Last modified
SRR_1.fastq	30 MB	Yesterday
SRR_2.fastq	30 MB	Yesterday

Files tab under DATA section showing uploaded files

Now that the samples are uploaded you will need to return and follow the steps in 3.0 to get your metadata associated with each sample and that should look like below.



Collection_date	iso_continent	iso_country	iso_county	iso_state
6/4/2021	North America	USA	Carson City	Nevada
6/10/2021	North America	USA	Clark County	Nevada

TABLE data showing metadata of the samples.

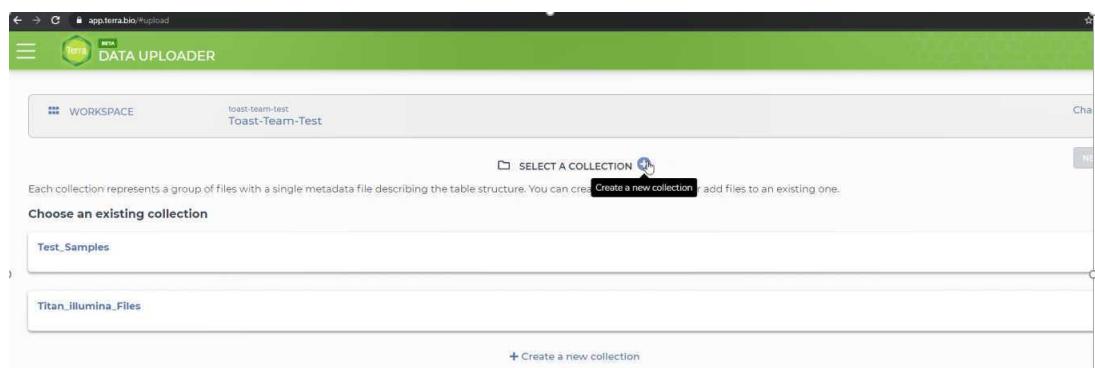
Here is a short video demonstration:

-
- 3.3 Alternatively, you can upload your files via '<https://app.terra.bio/#upload>'. There is no button on Terra to take you to this page, you will need to type this into the search bar.

Navigate to '<https://app.terra.bio/#upload>'

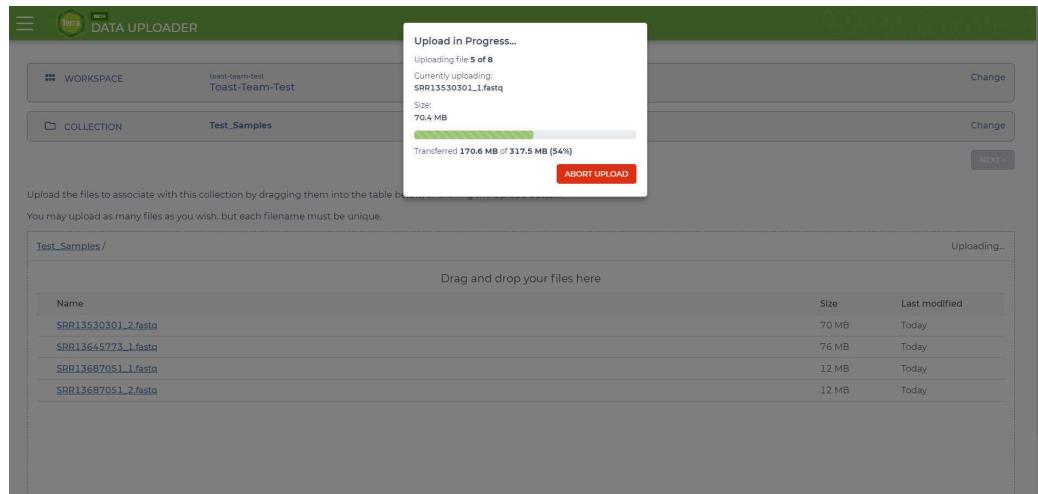
Select the Terra Workspace to which you would like to upload your fastq files. This will be the same workspace created in the previous step.

Click the '+ Create a new collection' link and enter a name for your new collection of fastq files. **DO NOT INCLUDE SPACES IN THE COLLECTION NAME, use underscores instead. Spaces will cause an error later in the pipeline.**



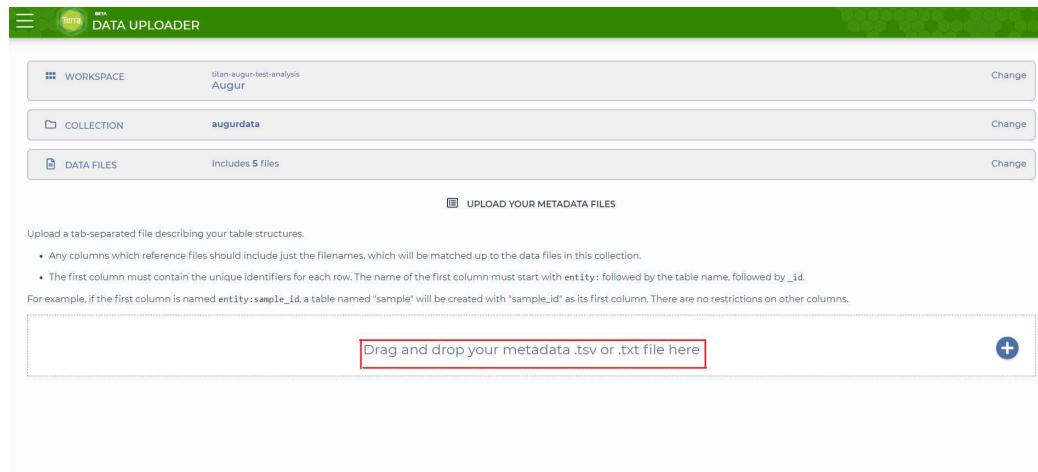
Click the 'Create Collection' button and you will be routed to the data uploader page where you can drag and drop fastq files directly into your browser window to initiate upload.

Drag and drop the fastq files that you would like to upload in to the upload space.



Screen when files are uploading

Once your files have been successfully uploaded, select "NEXT>" to proceed to the metadata upload page.



Drag and drop your table data

The screenshot shows the Terra Data Uploader interface. At the top, it displays the workspace 'titan-augur-test-analysis' and collection 'augurdata'. Below this, a table titled 'DATA FILES' shows 'Includes 5 files'. A button 'UPLOAD YOUR METADATA FILES' is present. In the center, a section titled 'Creating a new Table: titan_augur_sample' contains a table with the following data:

entity:titan_augur_sample_id	Read	Collection_date	Iso_country	Iso_state	Iso_continent
sample_1	SRR_1.fasta	6/4/2021	USA	Nevada	North America
sample_2	SRR_2.fasta	6/10/2021	USA	Nevada	North America
sample_3	SRR_3.fasta	6/10/2021	USA	NorthCarolina	North America
sample_4	SRR_4.fasta	6/12/2021	USA	Nevada	North America
sample_5	SRR_5.fasta	6/12/2021	USA	Ohio North America	Carson City

Buttons for 'CANCEL' and 'CREATE TABLE' are at the bottom right.

DATA UPLOADER page after creating new table metadata for the samples

Here is a video of the process:

Running the Augur Workflows

4 Running Titan_Augur is a two step process.

1. Run Titan_Augur_Prep workflow
2. Run Titan_Augur_Run workflow

4.1 Running Titan_Augur_Prep workflow

Navigate to the Augur Prep workflow page:

From 'DASHBOARD' menu -> select 'WORKFLOWS' -> select 'Titan_Augur_Prep' workflow tile

The screenshot shows the Terra WORKSPACES page. The 'WORKFLOWS' tab is selected. Below it, there is a grid of workflow tiles. One tile for 'Titan_Augur_Prep' is highlighted with a red box. The tile information is as follows:

- Name: Titan_Augur_Prep
- Version: V.v1.4.4
- Source: Dockstore
- Details icon (circled)

This will bring up the Titan_Augur_Prep workflow page:

Task name	Variable	Type	Attribute
titan_augur_prep	assembly	String	Required
titan_augur_prep	collection_date	String	Required
titan_augur_prep	iso_continent	String	Required
titan_augur_prep	iso_country	String	Required
titan_augur_prep	iso_state	String	Required
titan_augur_prep	pango_lineage	String	Required

Titan_Augur_Prep workflow page where you will import all the inputs from existing workspace

Select the version of the workflow you would like to run. If no preference pick the latest stable version (in this example v1.4.4). **The main and dev versions of the pipeline are under active development and are NOT recommended for users.**

Ensure that "Run workflow(s) with inputs defined by data table" is selected, in the Output tab select 'use_defaults', click the "Use call caching" box and then select the root entity type for the sample data you wish to analyze. **Don't use the "set" version!**

Call caching allows Terra to identify and skip jobs that have been run previously; this option is by default enabled to avoid unnecessary compute costs. More information on Terra call caching, including examples of when you may want to disable this feature, is available through the [Terra Support Documentation](#).

Click "SELECT DATA" and choose the samples you wish to analyze. **Note that only 25 samples are shown by default so make sure when you click the check box to select all (as shown in video) that you actually get all the samples you intend to!**

For the inputs to the workflow the top six rows in the input represent variables that have to be provided by the user. Complete the INPUTS section with the appropriate attributes. In our example, for the 'assembly_fastas' variable, for the 'Attribute' text box we set that to 'this.assembly.fasta' to indicate the 'fasta_files' we wish to analyze. The rest of the inputs are as follows:

- 'collection date' should be 'this.collection_date'
- 'iso_continent' should be 'this.iso_continent'
- 'iso_country' should be 'this.iso_country'
- 'iso_state' should be 'this.iso_state'
- 'pango_lineage' should be 'this.pango_lineage' --> like assembly fastas this obtained from the previous analysis (Titan or uploaded from your own).

Once you input all the data, it would look similar to the below picture.

The screenshot shows the 'Titan_Augur_Prep' workflow configuration page. At the top, there's a dropdown for 'Version' set to 'v1.4.4'. Below it, a 'Synopsis' section states 'No documentation provided'. Under 'Step 1', 'Run workflow(s) with inputs defined by data table' is selected. In 'Step 2', 'SELECT DATA' is chosen, and a note says '20 selected augur_prep_samples will create a new augur_prep_sample_set named "Titan_Augur_Prep_2021-07-01T17-13-35"'. There are checkboxes for 'Use call caching', 'Delete intermediate outputs', and 'Use reference disks'. Below these are tabs for 'SCRIPT', 'INPUTS' (highlighted with a red box), 'OUTPUTS' (highlighted with a red box), and 'RUN ANALYSIS'. A 'Hide optional inputs' link is available. On the right, there are 'Download Json' and 'SEARCH INPUTS' buttons. The main area displays a table of inputs:

Task name	Variable	Type	Attribute
titan_augur_prep	assembly	String	this.assembly.fasta
titan_augur_prep	collection_date	String	this.collection_date
titan_augur_prep	iso_continent	String	this.iso_continent
titan_augur_prep	iso_country	String	this.iso_country
titan_augur_prep	iso_state	String	this.iso_state
titan_augur_prep	pango_lineage	String	this.pango_lineage
prep_augur_metadata	CPU:s	Int	Optional

Once you are at this page and ready, hit Run and Launch the Analysis

NOTE: If you named your columns something other than fastas then just type "this." followed by whatever the column name is. We would advise naming your fastas column "fastas" for clarity.

Video showing on how to submit the Titan_Augur_Prep Run:

Once the Titan_Augur_Prep workflow runs successfully, you should be able to see the sample_metadata.tsv file associated to the sample table data under the TABLE section. This sample_metadata.tsv is selected as one of the inputs to run the Titan_Augur_Run workflow, that is why this step is necessary to check it was generated first.

The video below shows this step:

4.2

Running Titan_Augur_Run workflow

As we did before select the version of the workflow you would like to run (pick the latest stable version, here v1.4.4). **The main and dev versions of the pipeline are under active development and are NOT recommended for users.**

Ensure that "Run workflow(s) with inputs defined by data table" is selected and click the "Use call caching" box. **In contrast to the Titan_Augur_prep workflow, select "augur_prep_sample_set" in the the root entity section, because it will expect to have set of files in order to generate the output for the auspice**

visualizations. Since our root entity is a sample set when you click "SELECT DATA" click on the set of samples that you just ran Augur_Prep on.

And in INPUTS section, assembly_fastas attribute is set to "this.augur_prep_samples.assembly_fastas", Remember it will look for **set of samples** not just sample, so **name the attribute to root entity and set it to plural**.

The rest of the inputs are:

- Build_name in string format: Here we named it "B.1.147_samples_set", but you can call it whatever you like.
- sample_Metadata.tsvs (Generated by the Augur_Prep in the previous step) as "this.augur_prep_samples.augur_metadata. **Again note the plural form!**

If you don't use the plural form then you will get an error see section 8.2 for details.

And it should look like this.

Task name	Type	Attribute	Value
titan_augur_run	ONT_sample	this.augur_prep_samples.assembly_fastas	B1.147_all_samples_set
titan_augur_run	ONT_sample_set		
titan_augur_run	sample		
titan_augur_run	sample_set	d_name	B1.147_all_samples_set
titan_augur_run	sample_metadata_tsvs		
sarscov2_nextstrain	oncotools_trnns_to_infer		
sarscov2_nextstrain	auspice_config		
sarscov2_nextstrain	clades_tsv		

Once your input form is complete, move on to the OUTPUTS form and select "Use Defaults". Terra will then populate the OUTPUTS form with all of the default outputs options generated by the workflow. If you forget to do this you won't have easily accessible results! **Save these changes by clicking the 'Save' button**

Once your INPUTS and OUTPUTS forms are complete, click the 'Save' button on the top right-hand side of the page. The yellow caution icons should disappear and the Run Analysis option should be made available

You are now ready to run the Titan_Augur_Run workflow! Click on the 'Run Analysis' button to the right of the 'Outputs' tab. A popup window should appear titled 'Confirm launch'. If the 'Run Analysis' button is greyed out, you need to save your recent changes by clicking the 'Save' button.

Clicking the 'Launch' button should bring you to the 'Job History' panel where each sample will be queued for the Titan_Augur_Run analysis. The status will change from Queued -> Submitted -> Running

Video showing how to submit the Titan_Augur_Run workflow.

View and Download the Augur workflow output reports

- 5 First, **verify the sample set has successfully completed** by looking at the 'Workflow Status' section in the top left of the 'Job History' panel. The job has completed when the sample set have a status of 'succeeded' with a green checkmark.

Since Titan_Augur_Run analyzes a set of samples together, the job status would appear this way, it completed analysis successfully (hence the green checkmark).

Submission (click for details)	Data entity	No. of Workflows	Status	Actions	Submitted	Submission ID
Stan-evaluation_Titan_Augur_Run to dawenlong@gmail.com	Titan_Augur_Run_Prep_2021-07-11T13:38:33	1	✓ Done		Jul 8, 2021, 5:05 PM	1fb6e108-ae99-4bd6-92c5-9d70916abf00

Job Status

If you want to check what steps have finished within a job (helpful for troubleshooting), then click the job description under the submission column. This will take you to a page with more details on the job run.

Workflow Status:	Workflow Configuration	Submitted by	Total Run Cost:
✓ Succeeded: 1	Stan-evaluation_Titan_Augur_Run	dawenlong@gmail.com Jul 8, 2021, 5:05 PM	N/A
	Data Entity	Submission ID	Call Caching
	Titan_Augur_Run_Prep_2021-07-11T13:38:33	1fb6e108-ae99-4bd6-92c5-9d70916abf00	Enabled
	Workflow Details	Use Reference Disks	Disabled
	Delete Intermediate Outputs		
	Delete		

Workflow Details

Data Entity	Last Changed	Status	Run Cost	Messages	Workflow ID	Links
Titan_Augur_Run_Prep_2021-07-11T13:38:33 (augur_ares_sample_set)	Jul 8, 2021, 5:49 PM	✓ Succeeded	N/A	0	3fb4532c43e-462e-988e-f074-4235fce0	

Next click on the "Workflow Dashboard" icon on the far right under the links column. On the next page you will be able to tell what steps in the pipeline have successfully finished. Here all 3 tasks that were call finished successfully (hence the green checkmark).

Workflow Status	✓ Succeeded	Workflow Timing	Start: Jul 9, 2021, 10:33 AM
Calls		End:	Jul 9, 2021, 10:46 AM
Total Call Status Counts	✓ 3 Succeeded	Links	
Call Lists		Job Manager	Execution Directory
	> titan_augur_run.version_capture × 1		View execution log
	> titan_augur_run.sarscov2_nexstrain × 1		
	> titan_augur_run.snp_dists × 1		
	> Submitted workflow script		

Once confirming the job completed, go to the 'DATA' panel and under the 'TABLES' tab click on the sample table that you created and uploaded in step 4. It will be named 'sample (#)' which has the set of samples/entities in your file

The screenshot shows the Terra Data panel with the 'DATA' tab selected. Under the 'TABLES' section, there is a list of tables: 'ONT_sample' (20), 'ONT_sample_set' (2), 'assembly_sample' (9), 'assembly_sample_set' (1), 'augur_prep_sample' (20), 'augur_prep_sample_set' (1), and 'augur_sample' (1). The 'augur_prep_sample_set' table is highlighted with a red box. Below the table list is a preview of the data, showing columns like 'gur_prep_sample_set_id', 'augur_prep_samples', 'auspice_input_json', 'combined_assemblies', 'keep_list', and 'MAFFT'. A red box highlights the '20 entities' count.

Titan_Augur_Run output with added results after the analysis

Now the Terra sample table will have the additional attributes that were added by the workflow when you specified the output names (set to default in this example). You can reduce the number of fields you want to visualize by clicking the "gear" icon in top row on the right. Select only the fields you want to see then click "Done".

This screenshot is similar to the previous one, showing the Terra Data panel with the 'augur_prep_sample_set' table selected. A red box is placed over the gear icon in the top right corner of the table preview area, indicating where to click to manage column visibility.

The screenshot shows the Terra Data panel with the 'augur_prep_sample_set' table selected. A 'Select columns' modal window is open over the table preview. The modal lists various columns with checkboxes: 'augur_prep_samples' (checked), 'auspice_input_json' (checked), 'combined_assemblies' (checked), 'keep_list' (checked), 'MAFFT_alignment' (unchecked), 'metadata_merged' (unchecked), 'snp_matrix' (unchecked), 'time_tree' (unchecked), 'titan_augur_run_analysis_date' (unchecked), 'titan_augur_run_version' (unchecked), and 'unmasked_snps' (unchecked). At the bottom of the modal are 'CANCEL' and 'DONE' buttons.

Pop up menu with options of what metrics you want to see

The three files that are used for visualization are:

auspice_input_json - The output generated from the NextClade analysis step of the Augur workflow. This file includes the samples for clade typing and the single sample placed on the tree. Downloading this file won't be saved to your local folders but will be opened in a browser. Make sure

you right click on the page and "save as" it to your local directory

combined_assemblies - A concatenation of all of the assemblies that were included in this phylogenetic analysis. For this example it includes all of 20 samples combined into one single filtered fasta file combined_assemblies.

metadata_merged - Every metadata file from each sample that is generated by Titan_Augur_Prep is merged into one single metadata file.

The specific filenames produced by the pipeline can be viewed in the Terra sample report by scrolling to the right (or just only selecting those columns using the "gear" icon). You can download or copy this report by using either the 'Download All Rows' or 'Copy Page To Clipboard' buttons at the top of the table. To download a particular output file, click on the link under the column name a popup window titled 'File Details' should appear.

The screenshot shows the Terra Data interface. In the 'Tables' section, several datasets are listed. Three specific files are highlighted with red boxes: 'auspice_input_json', 'combined_assemblies', and 'metadata_merged'. The 'auspice_input_json' file is currently selected. The interface includes a 'DOWNLOAD ALL ROWS' button and a search bar.

The three primary output files of interest

Click the 'Download For < \$0.XX' button to save specific files to your local directory for further analyses or visualization.

The screenshot shows the Terra Data interface with a 'File Details' modal open. The modal displays the contents of the 'auspice-B1.147_all_samples_set_auspice.json' file, which is a JSON object containing metadata and assembly information. It shows the file size as 44.51 KB and provides a 'DOWNLOAD FOR < \$0.01*' button. Below the button, there is a terminal download command and a note about estimated download costs. The modal has a 'DONE' button at the bottom right.

Download all 3 files now to perform the phylogenetic visualization in the next step.

Video showing how to download and save the output files.

Auspice Visualization

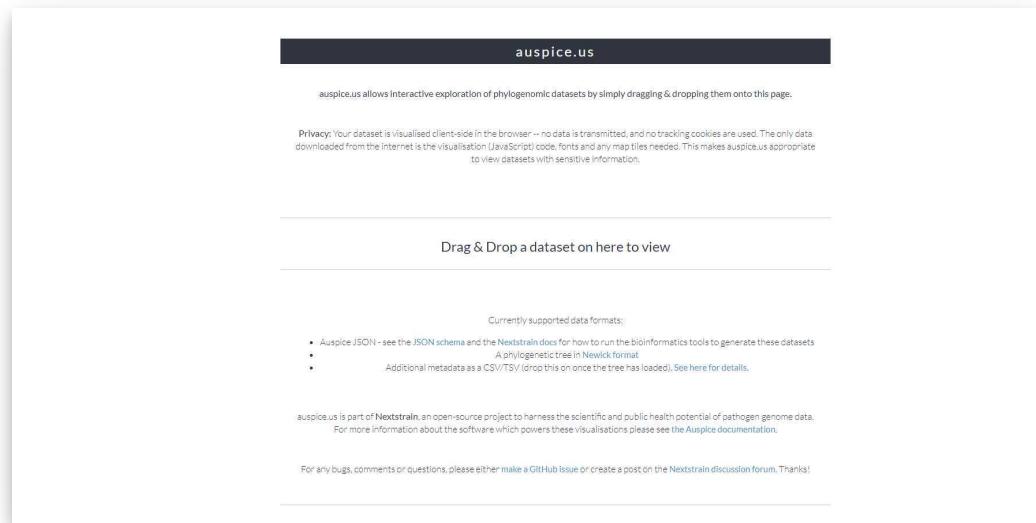
- 6 After completing the Augur workflow, the resulting output files contain the comparative metrics necessary to analyze the genetic relatedness between samples and draw phylogenetic inference by visualization with Auspice.

Auspice takes the json file to build an interactive phylogenetic tree with Nextstrain. Visit the [CDC COVID-19 Genomic Epidemiology Toolkit](#) for an introduction to these applications:

Module 3.2: Getting started with Nextstrain

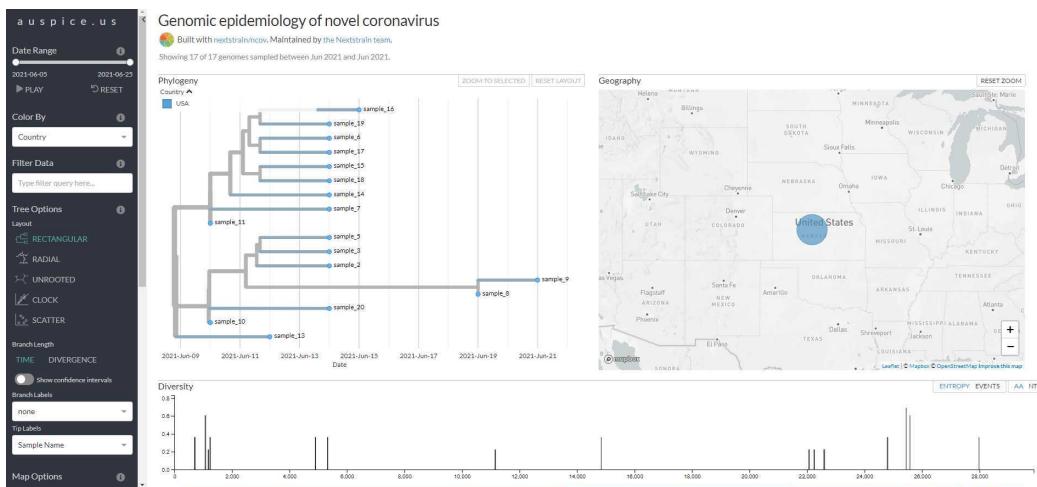
Module 3.4: Walking through Nextstrain trees

You should have downloaded auspice_input_json and metadata merged files from Step 5. We will now use them to generate a phylogenetic visualization using [auspice.us](#).



Auspice.us home page

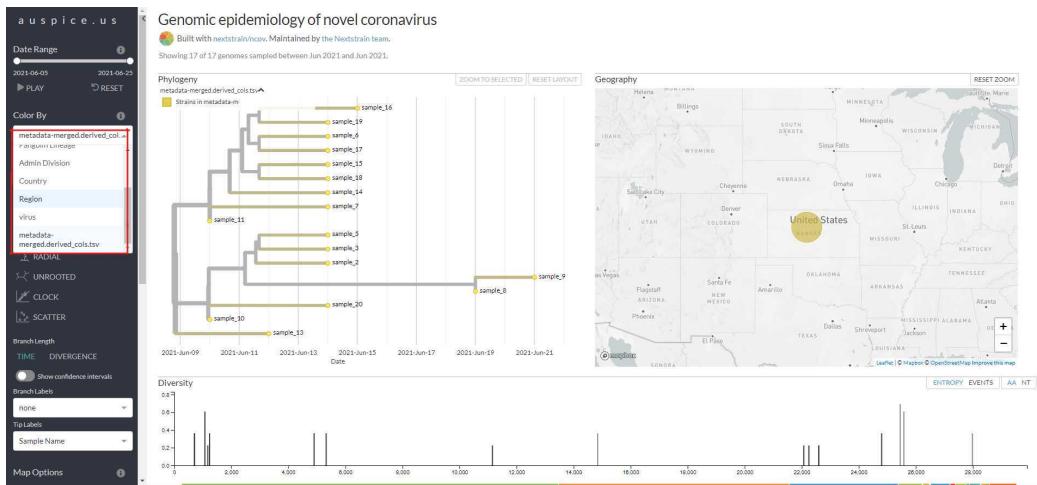
Data can be uploaded to Auspice by simply drag-and-drop of the auspice_input_json file, which should automatically visualize a phylogenetic tree similar the example below:



Example Auspice visualization

Hover over the tree and geographical map to see the divergence of the strains, toggle between the options and parameters based on your needs.

Metadata can be added to annotate the tree for assistance with identify local outbreaks or clustered samples by simply drag-and-drop of the metadata_merged file onto the tree.



After adding metadata

Note: Once you drop the metadata file, all the attributes will show up in the "Color By" section (like in the picture above) where you can visualize your samples by whatever attribute that you have in your metadata file.

Visualization gives few default errors and it would not be a problem, but most importantly it gives the option to color code your samples by metadata variables.

Since the metadata in this example is random dummy data so results are not accurate, but if you upload your metadata, it would allow you to 'Color By' the metadata variables.

This video shows how to drop and visualize on auspice.

UShER Visualization

7 It may also be helpful to visualize your results using UCSC's UShER webportal.

Visit the [CDC COVID-19 Genomic Epidemiology Toolkit](#) for an introduction to UShER:

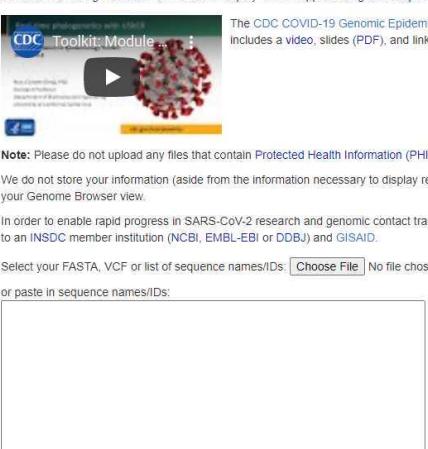
Module 3.3: Real-time phylogenetics with UShER

Start by navigating to the UShER SARS-CoV-2 webportal: <https://genome.ucsc.edu/cgi-bin/hgPhyloPlace>

UShER: Ultrafast Sample placement on Existing tRee

Upload your SARS-CoV-2 sequence (FASTA or VCF file) to find the most similar complete, high-coverage samples from GISAID or from public sequence databases (NCBI Virus / GenBank, COG-UK and the China National Center for Bioinformation), and your sequence's placement in the phylogenetic tree generated by the sarscov2phylo pipeline. Placement is performed by Ultrafast Sample placement on Existing tRee (UShER) (Turakhia et al.). UShER also generates local subtrees to show samples in the context of the most closely related sequences. The subtrees can be visualized as Genome Browser custom tracks and/or using Nextstrain's interactive display which supports drag-and-drop of local metadata that remains on your computer.

The CDC COVID-19 Genomic Epidemiology Toolkit now includes a training module for UShER Module 3.3 includes a video, slides (PDF), and links to more resources.



Note: Please do not upload any files that contain Protected Health Information (PHI) to UCSC.
We do not store your information (aside from the information necessary to display results) and will not share it with others unless you choose to share your Genome Browser view.
In order to enable rapid progress in SARS-CoV-2 research and genomic contact tracing, please share your SARS-CoV-2 sequences by submitting them to an INSDC member institution (NCBI, EMBL-EBI or DDBJ) and GISAID.

Select your FASTA, VCF or list of sequence names/IDs: No file chosen
or paste in sequence names/IDs:

UShER SARS-CoV-2 home page

Upload your concatenated fasta file (combined_assembly from Step 5) by clicking "Choose File" and selecting the combined_assembly file then click "Upload". You can also open the combined_assembly file and then copy and paste the sequence names/IDs.

UShER is very fast and will place your sequences within a phylogenetic tree of public sequences in [GISAID](#) or GenBank within few minutes. After uploading your assembly file pick which phylogenetic tree version you want to use from the drop down menu. **These are updated over time as new genomes are added so its a good idea to take note of the version you choose to use.**

UShER outputs subtree results with the 50 sequences (default) most-closely related to your input samples.

You can pick the Phylogenetic tree version if you want to see the from where/how your samples are diverged, UShER would reproduce the results adding your samples to the existing tree nodes. We

recommend using the top, most up to date tree.

Subtrees generated by UShER depend on the divergence among the input samples, if the strains are more divergent, it will produce more subtrees, and if you wish to add more samples, UShER can add them to the existing tree based on sample divergence and mutation. Trees can then be visualize and explore this using the Nextstrain platform.

In addition to phylogenetic placement, UShER also provides extra information, including QC metrics, in table format on the UShER results page.

If you wish to view the subtrees, select the options that is available below the table to launch the phylogenetic tree on the Nextstrain platform. Metadata can be added by drag-and-drop similar to with Auspice (Step 6).

For this demo UShER took the concatenated assembly fastas of 20 samples and constructed a phylogenetic tree against all GISAID trees that are being maintained by the UShER and it gave 16 subtrees (See "Subtree Number" column).

If you hover over the "?" of the column variables/sample rows, a gray dialogue box would appear (like in the picture) showing the info on that attribute.

	Downloads: Global phylogenetic tree with your sequences TSV summary of sequences and placements TSV summary of Spike mutations ZIP file of subtrees JSON and Newick files															
	Fasta Sequence	Size (?)	#Ns (?)	#Mixed (?)	Bases aligned (?)	Inserted bases (?)	Deleted bases (?)	#SNVs used for placement (?)	#Masked SNVs (?)	Nextstrain clade (?)	Neighboring sample in tree (?)	Lineage of neighbor (?)	#Imputed values for mixed bases (?)	#Maximally parsimonious placements (?)	Parsimony score (?)	Subtree number (?)
sample_10	29782 (?)	314	0	29782 (?)	0	0	4 (?)	0	20C	WalesPHC-2996F/2020 2020-04-03	?	0	50	0	2 (view in Nextstrain)	
sample_11	29782 (?)	0 (?)	0	29782 (?)	0	0	6 (?)	0	20A	IHUCOVID-0483 LR794668.1 ?	B.1.147	0	1	0	3 (view in Nextstrain)	
sample_13	29773 (?)	0 (?)	0	29773 (?)	0	9 (?)	9 (?)	0	20A	USA/AM-MD-HS-SC2014B/2020 MT439273.1 2020-03-24	B.1	0	1	3	4 (view in Nextstrain)	
sample_14	29782 (?)	0 (?)	0	29782 (?)	0	0	8 (?)	0	20A	IHUCOVID-0483 LR794668.1 ?	B.1.147	0	1	2	5 (view in Nextstrain)	
sample_15	29782 (?)	1649	0	27916 (?)	0	bases 55 - 29836 align to reference bases 55 - 29836	6 (?)	0	20A	IHUCOVID-0483 LR794668.1 ?	B.1.147	0	2	0	6 (view in Nextstrain)	
sample_16	29782 (?)	0 (?)	0	29782 (?)	0	0	6 (?)	0	20A	IHUCOVID-0483 LR794668.1 ?	B.1.147	0	1	0	7 (view in Nextstrain)	
sample_17	29782 (?)	0 (?)	0	29782 (?)	0	0	6 (?)	0	20A	IHUCOVID-0483 LR794668.1 ?	B.1.147	0	1	0	8 (view in Nextstrain)	
sample_18	29454 (?)	1040	0	28406 (?)	0	0	5 (?)	0	20A	IHUCOVID-0483 LR794668.1 ?	B.1.147	0	3	0	9 (view in Nextstrain)	
sample_19	29782 (?)	526	0	29259 (?)	0	0	6 (?)	0	20A	IHUCOVID-0483 LR794668.1 ?	B.1.147	0	2	0	10 (view in Nextstrain)	
sample_2	29775 (?)	0 (?)	0	29775 (?)	0	9 (?)	7 (?)	0	20C	USA/UT-UPLH-20121082/2020 MV210972.1 2020-06-19	B.1.146	0	26	1	11 (view in Nextstrain)	
sample_20	29782 (?)	2356	0	27426 (?)	0	0	5 (?)	0	20A	England/UK-9E869/2020 2020-04-11	?	0	507	1	12 (view in Nextstrain)	
sample_3	29782 (?)	2144	0	27638 (?)	0	0	6 (?)	0	20C	USA/UT-UPLH-20121082/2020 MV210972.1 2020-06-19	B.1.146	0	235	0	13 (view in Nextstrain)	
sample_5	29782 (?)	2144	0	27638 (?)	0	0	6 (?)	0	20C	USA/UT-UPLH-20121082/2020 MV210972.1 2020-06-19	B.1.146	0	235	0	14 (view in Nextstrain)	
sample_6	29782 (?)	0 (?)	0	29782 (?)	0	0	6 (?)	0	20A	IHUCOVID-0483 LR794668.1 ?	B.1.147	0	1	0	15 (view in Nextstrain)	
sample_7	29782 (?)	n (?)	n	29782 (?)	n	n	9 (?)	n	20A	IHUCOVID-0483 LR794668.1 ?	B.1.147	n	1	5	16 (view in	

Output metrics generated by UShER tool

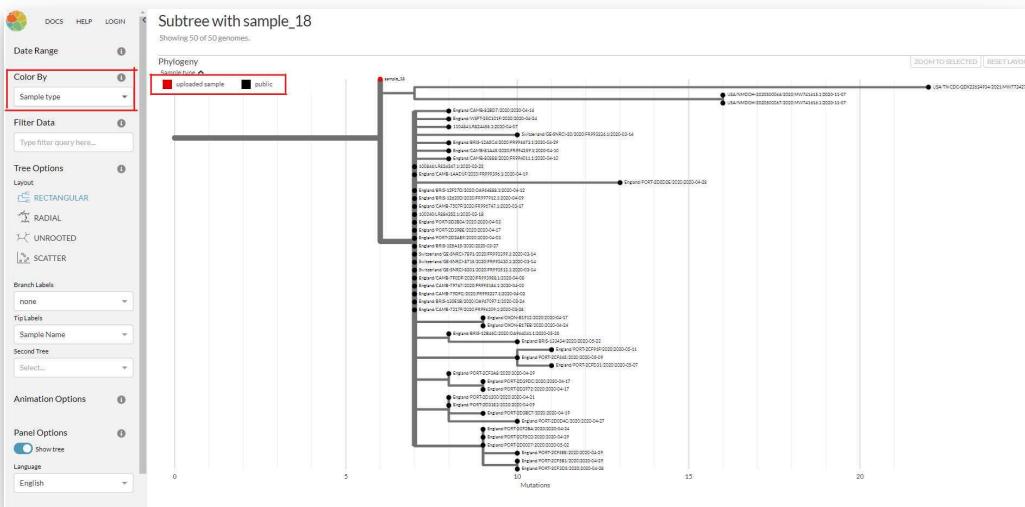
Colors are represented to show the variations of each attribute between samples.

The lower the score the lighter the color, which means you samples are of high quality when it is green. When the score is high the output would appear in Red. In general the higher the score/red in color in the metrics, suggests that the sample sequence has many errors and sample placements might not be reliable.

As an example: In the image above the Parsimony score, which explains the number of mutations on the branch, for Sample_10 is highlighted in red, this says the sample has relatively more mutations added to its branch than the others.

If you want to see the subtrees click the subtree section of the row you are interested in and the link associated to that sample will redirect to the subtree visual rendered on the Nextstrain platform. Nextstrain distinguishes the public samples in Black and your input samples in Red. Subtrees basically shows how your sample is diverged and placed among the public sequences.

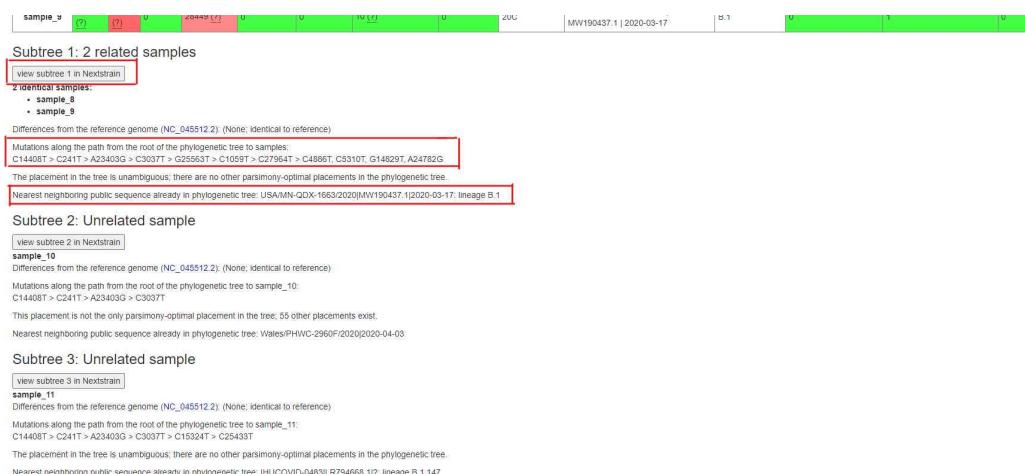
Here is an example subtree:



To view the subtree, sample_18 was selected and opted to color by sample type from the menu on left side. Nextstrain distinguishes the public samples in Black and your input samples in Red.

You can also see the subtrees by scrolling down the page and selecting from the individual subtree.

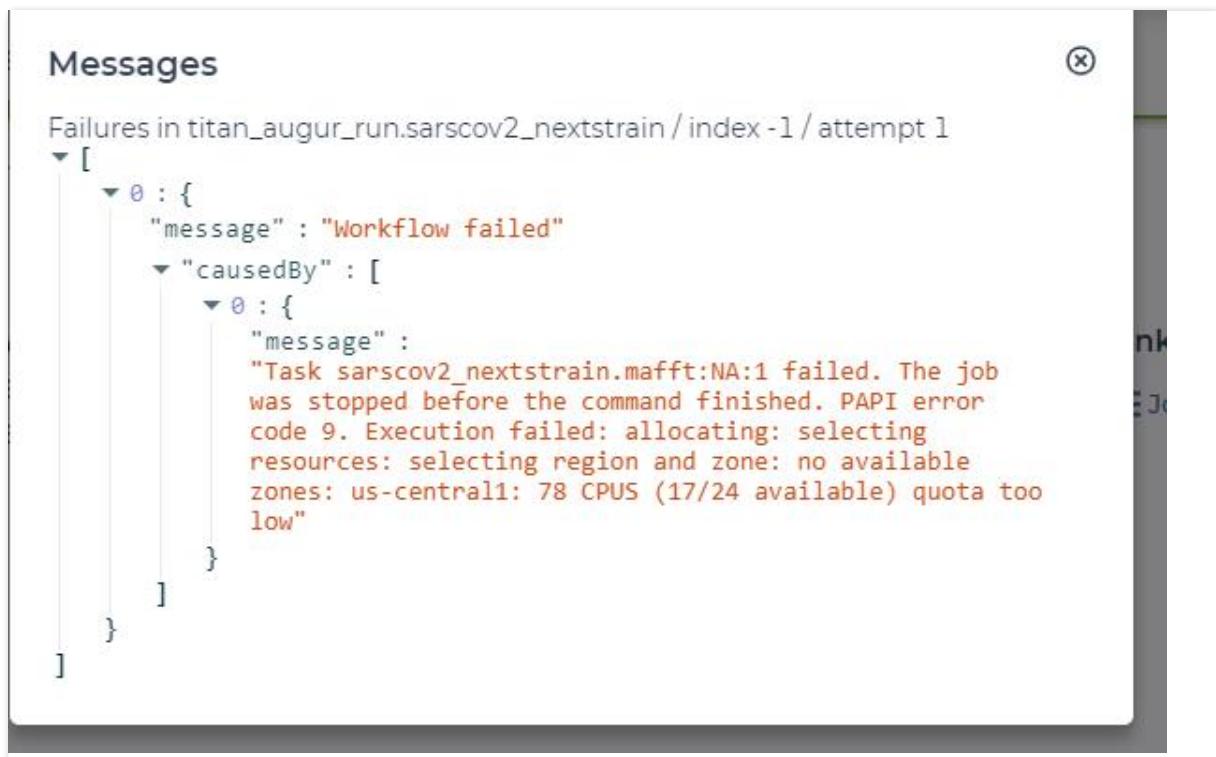
Additional information, such as the likely path of the mutation from root of the tree to your sample, most closely-related sequence in the global tree, and PANGO lineage inference, are provided.



How to debug a failed run

1. Failure due to the unavailability of the compute environment.

The Titan_Augur_Run workflow script requests 78 CPUs in order to perform the phylogenetic reconstruction for visualization in Auspice and Nextstrain. If the analysis fails due to lack of resources, try requesting for more CPUs to Terra.



The screenshot shows a 'Messages' panel with the title 'Messages' and a close button. Below the title, it says 'Failures in titan_augur_run.sarscov2_nextstrain / index -1 / attempt 1'. A dropdown arrow reveals a JSON object:

```
[{"0": {"message": "Workflow failed", "causedBy": [{"0": {"message": "Task sarscov2_nextstrain.mafft:NA:1 failed. The job was stopped before the command finished. PAPI error code 9. Execution failed: allocating: selecting resources: selecting region and zone: no available zones: us-central1: 78 CPUS (17/24 available) quota too low"}]}]}
```

For more details on how to fix this see [step 2 of our Terra Troubleshooting protocol](#).

2. Failure for not specifying the input attributes correctly.

There are two flavors of this error for Titan_Augur_Run:

- Singular instead of plural notation
- Using sample instead of Sample_Set

First, the singular instead of plural problem:

The Titan_Augur_run workflow looks for concatenated set of fasta files, **you need to set the root entity to "name_sample_set" (with name being specific to your dataset) and the input attribute notation of "this.name_samples.assembly.fasta and metadata as "this.name_samples.metadata"** rather than using "this.assembly.fasta" and "this.metadata" as we did with the Augur_prep protocol.

In other words, **you select the name_sample_set (singular) in root entity, however for the input use name_samples (plural) notation**, since it is looking for all the files as concatenated set.

You will see this error if you didn't use the plural:

Dashboard Data Notebooks Workflows Job History

↳ Job History > Submission 5d81c5bd-d32d-41c1-9e51-4c7aaa74eb19 > Workflow 1ec376a8-d348-4755-a80e-d34a5c962275
Workflow metadata fetched in 423ms

Workflow Status

⚠ Failed

Workflow Timing

Start: Jul 8, 2021, 5:02 PM
End: Jul 8, 2021, 5:02 PM

Links

Job Manager Execution Directory View execution log

↳ Workflow-Level Failures

↳ [↳ { ↳ { ↳ { "message": "Workflow input processing failed", "causedBy": [↳ { ↳ { "message": "Failed to evaluate input 'assembly_fastas' (reason 1 of 1): An Array[File]+ must contain at least one element" }, ↳ { ↳ { "message": "Failed to evaluate input 'sample_metadata_tsvs' (reason 1 of 1): An Array[File]+ must contain at least one element" }] }] }] }

↳ Calls

Total Call Status Counts
No calls have been started by this workflow.

↳ Submitted workflow script

Example error caused by singular instead of plural in inputs

Second, we will look at the problem of using sample instead of sample set as your root entity type:

Dashboard Data Notebooks Workflows Job History

Workflow History > Submission bcf02d924-8cda-4dfe-d220-2e0e2f0505d5 > Workflow 53843702-9641-4487-a505-d340a7b40000

Workflow metadata finished in 424ms.

Workflow Status

⚠ Failed

Workflow-Level Failures

- Workflow failed
 - message: "Workflow failed"
 - causedBy:
 - message: "Workflow failed"

Call Lists

Total Call Status Counts

- ✓ Success
- ⚠ Failed

Call Lists

> title_augr_rsn_veravia_capture + 1

✓ title_augr_rsn_sarscov2_hextrain + 1

Index Attempt Status Start End Call Caching Result Links

N/A	1	⚠ Failed	Jul 8, 2021, 5:59 PM	Jul 8, 2021, 4:15 PM		⚠ 1 Message
-----	---	----------	----------------------	----------------------	--	-------------

Submitted workflow script:

Error in Terra.bio job history

The specific error message can be found in the `draft_augur_tree.log` file in the specified Google Cloud storage location.

Example draft augur tree.log output with error message

To avoid this error, ensure the root entity type is set to 'sample_set' and not 'sample', as shown below.

The screenshot shows the Augur UI interface for a workflow named 'Titan_Augur_Run'. The 'WORKFLOWS' tab is selected. A specific step, 'Step 2', is highlighted. The step details are as follows:

- Step 1**: 'augur_prep_sample' (highlighted with a red box)
- Step 2**: 'SELECT DATA' (highlighted with a red box)

The 'SELECT DATA' step has the following configuration:

- Entity type: 'augur_prep_sample'
- Use call caching: checked
- Script: 'augur_prep_sample' (highlighted with a red box)
- Inputs: 'augur_prep_sample' (highlighted with a red box)
- Outputs: 'augur_prep_sample' (highlighted with a red box)

The 'augur_prep_sample' input table is displayed below:

Task name	Value	Type	Attribute
titian_augur_run	ONT_sample	file	this.augur_prep_samples.assembly.fasta
titian_augur_run	ONT_sample.set	array[files]	
titian_augur_run	sample.fasta	array[files]	
titian_augur_run	sample.set	array[files]	
titian_augur_run	sample_metadata.csv	array[files]	this.augur_prep_samples.augur_metadata
sanscov2_neextstrain	ancestral_muts_all_lnter	array[string]	Optional
sanscov2_neextstrain	basepair.config	file	Optional
sanscov2_neextstrain	clade.csv	file	Optional
sanscov2_neextstrain	clock_rate	float	Optional
sanscov2_neextstrain	clade_mtdev	float	Optional
sanscov2_neextstrain	rot_longs.tsv	file	Optional
sanscov2_neextstrain	min_unambig_genome	int	Optional
sanscov2_neextstrain	ref.fasta	file	Optional

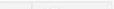
At the bottom right, there are buttons for 'Download json' and 'SEARCH INPUTS'.

Make sure the root entity type is 'sample_set'

3. Error due to spacing in collection name

In this example of a failed run we have an error caused by adding spaces into your collection name. If there is a failure for some other reason you can follow similar steps to understand why your job failed. Here we will be using Illumina paired-end data, but the process of debugging an error is the same if you have ONT data.

If a run fails you will see this indicated in the job history screen in the "status" column.

Job History		Submitted by		Total Run Cost			
Workflow Statuses		Workflow Configuration		Call Caching			
⚠ Failed: 1		toast-team-test/Ilan_illumina_PE		Disabled			
Data Entity 3015780998_ZZYGIWY sample		Submission ID 97f5da29-71d1-45e8-9938-82ba0a36b788		Disabled			
Delete Intermediate Outputs Disabled		User Reference Disk Disabled					
Search	Completion status	Last Changed	Status	Run Cost	Messages	Workflow ID	Links
3015780998_ZZYGIWY [sample]	May 13, 2021, 5:46 PM	⚠ Failed	N/A		e1687504-3421-4a20-9de1-11e3542233e3		

Job Failure

To understand why it failed click on the "workflow dashboard" icon in the "links" column. This will take you to a new screen and you can click the arrows next to the "message" to expand the message and see what it says. Here we see there are two errors that direct us to a log file to check. To find out more click on the "execution directory" icon under the links header. This will take you to the google bucket with all the output from the run.

WORKSPACES Workspaces - Local Learn Test Job Failure Example : Job History

DASHBOARD DATA NOTEBOOKS WORKFLOWS JOB HISTORY

← Job History · Submission 97f5da29-71d1-45e8-993b-82ba0a36b708 · Workflow e1687504-3421-4a20-9de1-11e35422d3e3 · Workflow metadata fetched in 446ms

Workflow Status Failed

Workflow-Level Failures □

Workflow Timing Start: May 13, 2021, 5:41 PM End: May 13, 2021, 5:45 PM

Links Job Manager Execution Directory View execution log Execution directory

Total Call Status Counts

Job failure messages

Follow the file path to the Google Cloud storage bucket location to the log files referenced in the error messages (in red/orange text in the above photo).

Cloud Storage							Bucket details		REFRESH		LEARN	
OBJECTS		CONFIGURATION		PERMISSIONS		RETENTION		LIFECYCLE				
Buckets	>	fc-650035cc-5856-433b-95e6-27c1f5cfcb7e										
UPLOAD FILES	UPLOAD FOLDER	CREATE FOLDER	MANAGE HOLDS	DOWNLOAD	DELETE							
Filter by name prefix only ▾												
	Name	Size	Type	Created time	Storage class	Last modified	Public access	Encryption	Retention expiration date	Holds		
<input checked="" type="checkbox"/>	fastlog/_raw.log	4 kB	text/plain; charset=UTF-8	May 13, 2021, 9:40:11 PM	Standard	May 13, 2021, 9:40:11 PM	Not authorized	Google-managed key	—	None		
<input checked="" type="checkbox"/>	gcs_debounce.sh	4.5 kB	text/plain; charset=UTF-8	May 13, 2021, 9:41:29 PM	Standard	May 13, 2021, 9:41:29 PM	Not authorized	Google-managed key	—	None		
<input checked="" type="checkbox"/>	gcs_localization.sh	1.6 kB	text/plain; charset=UTF-8	May 13, 2021, 9:41:29 PM	Standard	May 13, 2021, 9:41:29 PM	Not authorized	Google-managed key	—	None		
<input checked="" type="checkbox"/>	gcs_transfer.sh	13.4 kB	text/plain; charset=UTF-8	May 13, 2021, 9:41:29 PM	Standard	May 13, 2021, 9:41:29 PM	Not authorized	Google-managed key	—	None		
<input checked="" type="checkbox"/>	pipelines/logo/	—	Folder	—	—	—	—	—	—	—		
<input checked="" type="checkbox"/>	rc	2 kB	text/plain; charset=UTF-8	May 13, 2021, 9:42:26 PM	Standard	May 13, 2021, 9:42:26 PM	Not authorized	Google-managed key	—	None		
<input checked="" type="checkbox"/>	script	1.9 kB	text/plain; charset=UTF-8	May 13, 2021, 9:41:29 PM	Standard	May 13, 2021, 9:41:29 PM	Not authorized	Google-managed key	—	None		
<input checked="" type="checkbox"/>	stderr	582 kB	text/plain; charset=UTF-8	May 13, 2021, 9:42:23 PM	Standard	May 13, 2021, 9:42:23 PM	Not authorized	Google-managed key	—	None		
<input checked="" type="checkbox"/>	stdout	49 kB	text/plain; charset=UTF-8	May 13, 2021, 9:42:24 PM	Standard	May 13, 2021, 9:42:24 PM	Not authorized	Google-managed key	—	None		

Click on log file referenced in the error message.

Click on the "Authenticated URL" link that will take you to a text file.

Cloud Storage

Object details

Download Edit metadata Edit permissions Delete

Browser Monitoring Settings

Since you are not authorized to know the public access status of this object, it is possible that the public URL displayed is not valid.

Overview

Type	text/plain; charset=UTF-8
Size	4 kB
Created	May 13, 2021, 5:45:11 PM
Last modified	May 13, 2021, 5:45:11 PM
Custom time	—
Public URL	https://storage.cloud.google.com/fc.a40039e0-98d5-43f8-b965-7517fcfb7a/g77fda70-7181-44ab-9403-034b02b7b688?Expires=1621447340&Signature=1e7a42d21a3c0a840323337378&HttpMethod=GET&CacheControl=max-age%3D0
Authenticated URL	https://storage.cloud.google.com/fc.a40039e0-98d5-43f8-b965-7517fcfb7a/g77fda70-7181-44ab-9403-034b02b7b688?Expires=1621447340&Signature=1e7a42d21a3c0a840323337378&HttpMethod=GET&CacheControl=max-age%3D0
gsutil URL	gs://fc.a40039e0-98d5-43f8-b965-7517fcfb7a/g77fda70-7181-44ab-9403-034b02b7b688?Expires=1621447340&Signature=1e7a42d21a3c0a840323337378&HttpMethod=GET&CacheControl=max-age%3D0

Permissions

Public access	Not authorized
Protection	None
Hold status	None
Retention policy	None

Authenticated URL link

In the text file we can see that there was an error that cause by there being a space between "Bad"

and "Sample" in our file path that was created when we made "Bad Sample" rather than "Bad_Sample" as our collection name.

Example fastqc_raw.log file showing the error

[Video of the whole process.](#)

Some error message can be difficult to understand. Contact **TOAST@cdc.gov** for help debugging job failures.