

<b>МІЖНАРОДНА СТАНДАРТНА ОПЕРАЦІЙНА ПРОЦЕДУРА PULSENET ДЛЯ АНАЛІЗУ ДАНИХ КОРОТКОГО ЗЧИТУВАННЯ WGS ILLUMINA З ВИКОРИСТАННЯМ ПЛАТФОРМИ TERRA.BI</b>			
Док. No. PNID01	Вер. No. 01	Дата набуття чинності:	Сторінка 1 з 67

## **ЗМІСТ - ГІПЕРПОСИЛАННЯ НА ПРОЦЕДУРИ**

- [Підготуйте файл метаданих для послідовностей, які будуть завантажені до Terra \(5.1\)](#)
- [Увійдіть на Terra за допомогою Chrome та свого облікового запису Google \(5.2\)](#)
- [Завантажте файли послідовностей та метадані до Terra \(5.3\)](#)
- [Запустіть робочий процес контролю якості та генотипування \(5.4\)](#)
- [Оцініть метрики контролю якості для послідовностей \(5.5\)](#)
- [Перегляньте результати генотипування для послідовностей \(5.6\)](#)
- [Завантаження послідовностей до NCBI \(5.7\)](#)
- [Додаток PNID01-1: Імпорт даних до Terra безпосередньо з Illumina BaseSpace](#)
- [Додаток PNID01-2: Завантаження даних з NCBI SRA](#)
- [Додаток PNID01-3: Налаштування подання таблиці даних для показників контролю якості PulseNet](#)
- [Додаток PNID01-4a. Критичні показники якості PulseNet для подачі рутинних послідовностей](#)
- [Додаток PNID01-4b. Етап попереднього скринінгу зчитування TheiaProk для виключення неякісних послідовностей з метою економії обчислювальних ресурсів](#)
- [Додаток PNID01-5. Налаштування подання таблиці даних для генотипування аналізів PulseNet](#)
- [Додаток PNID01-6. Завантаження додаткових метаданих до Terra для подання до NCBI та налаштування вигляду таблиці даних для метаданих](#)
- [Додаток PNID01-7: Огляд робочого процесу TheiaProk для визначення бактеріальних характеристик](#)

1. **МЕТА:** Описати процедуру аналізу даних повногеномного секвенування (WGS) з коротким зчитуванням Illumina, які будуть використані для спостереження PulseNet International (PNI) з використанням хмарної платформи Terra.Bio.
2. **СФЕРА ЗАСТОСУВАННЯ:** Дана процедура поширюється на весь персонал PulseNet, який використовує платформу Terra.Bio для аналізу даних секвенатора Illumina для короткого зчитування WGS в рамках діяльності з нагляду в міжнародній мережі PulseNet. Дана СОП охоплює завантаження послідовностей і метаданих до Terra.Bio, оцінку якості послідовностей, робочі процеси збірки і генотипування, а також завантаження послідовностей до NCBI. Філогенетичний аналіз охоплюється СОП PNID02 (PulseNet International Standard Operating Procedure for Phylogenetic Analysis of WGS Data Using the Terra.Bio Platform - Стандартна операційна процедура для філогенетичного аналізу даних WGS з використанням платформи Terra.Bio).
3. **ВИЗНАЧЕННЯ/ТЕРМІНИ:**
  - 3.1 **ANI:** Середня нуклеотидна ідентичність.
  - 3.2 **BaseSpace:** Хмарне обчислювальне середовище Illumina для аналізу, управління та зберігання даних секвенування наступного покоління, включаючи обмін даними.

<b>МІЖНАРОДНА СТАНДАРТНА ОПЕРАЦІЙНА ПРОЦЕДУРА PULSENET ДЛЯ АНАЛІЗУ ДАНИХ КОРОТКОГО ЗЧИТУВАННЯ WGS ILLUMINA З ВИКОРИСТАННЯМ ПЛАТФОРМИ TERRA.BI</b>			
Док. No. PNID01	Вер. No. 01	Дата набуття чинності:	Сторінка 2 з 67

- 3.3 Команди Bash:** Bash (Bourne Again Shell) - це оболонка інтерфейсу командного рядка (CLI), яка широко використовується в Linux і macOS. Оболонка - це комп'ютерна програма, яка дозволяє безпосередньо керувати операційною системою комп'ютера. Команди Bash використовуються для керування комп'ютером або операційною системою без необхідності навігації по меню, опціях і вікнах у графічному інтерфейсі користувача.
- 3.4 Біопроект:** Колекція біологічних даних про NCBI, пов'язаних з однією ініціативою, що походить від однієї організації або консорціуму.
- 3.5 Біозразки:** Надана описова інформація (метадані) про біологічні матеріали, з яких отримані дані, що зберігаються в NCBI.
- 3.6 Contig:** Суцільна консенсусна послідовність, отримана в результаті складання багатьох коротких фрагментів ДНК, що перекриваються.
- 3.7 Покриття:** Середня кількість зчитувань, які включають певний нуклеотид у реконструйовану послідовність.
- 3.8 Критичні показники якості:** покриття (після тримісії), середня якість (показник Q до тримісії), довжина збірки та чисельність вторинних родів (виявлення забруднення за допомогою MIDAS). Послідовності, які не відповідають мінімальним порогам/прийнятним діапазнам для цих метрик, визначених у цьому документі, слід повторно секвенувати.
- 3.9 CSV:** Значення, розділені комами.
- 3.10 Збірка DeNovo:** Збірка послідовностей, створена з коротких необроблених зчитувань без використання референтного геному.
- 3.11 FASTA:** текстовий формат для представлення послідовностей нуклеотидів або пептидів, в якому пари основ або амінокислот представлені за допомогою однолітерних кодів. Послідовність у форматі FASTA починається з однорядкового опису, за яким слідує рядки даних послідовності. Рядок опису відрізняється від даних послідовності символом "більше" (">") у першому стовпчику.
- 3.12 FASTQ:** текстовий формат для зберігання як біологічної послідовності, так і відповідних оцінок якості.
- 3.13 GAMBIT:** Метод геномної апроксимації для ідентифікації та відстеження бактерій. Метод ідентифікації бактеріальних видів, що використовує алгоритм на основі k-мер для пошуку у великій довідковій базі даних геномів.
- 3.14 Gzip:** Формат файлів і програмне забезпечення, що використовується для стиснення та розпакування файлів для швидкої передачі даних через Інтернет.
- 3.15 LIMS:** Система управління лабораторною інформацією.
- 3.16 Mash Sketching:** Mash - це набір інструментів для створення і використання скетчів MinHash, спосіб перетворення геному в невеликий підпис, який можна легко порівняти з іншими підписами.
- 3.17 Метадані:** Набір даних, які описують та надають інформацію про інші дані.
- 3.18 MIDAS:** Метагеномна система аналізу внутрішньовидового різноманіття. Інтегрований обчислювальний конвеєр для кількісної оцінки чисельності бактеріальних видів та покриття метагеномів на основі бласт-вирівнювання з панеллю універсальних однокопійних генів.

<b>МІЖНАРОДНА СТАНДАРТНА ОПЕРАЦІЙНА ПРОЦЕДУРА PULSENET ДЛЯ АНАЛІЗУ ДАНИХ КОРОТКОГО ЗЧИТУВАННЯ WGS ILLUMINA З ВИКОРИСТАННЯМ ПЛАТФОРМИ TERRA.BI</b>			
Док. No. PNID01	Вер. No. 01	Дата набуття чинності:	Сторінка 3 з 67

- 3.19 N50:** Статистика N50 зазвичай використовується для грубої оцінки геномних збірок. Вона показує довжину контига (в базових парах), для якої половина послідовності геному зібрана в контиги, розмір яких більший або дорівнює розміру контига N50.
- 3.20 NCBI:** Національний центр біотехнологічної інформації.
- 3.21 PNI:** PulseNet International.
- 3.22 QA/QC:** Забезпечення якості/контроль якості.
- 3.23 Оцінка якості:** Оцінка якості для кожної окремої позиції бази в послідовності, що вказує на точність розпізнавання нуклеотидних основ. Використовуються оцінки Phred, де  $Q = -10\log$  (ймовірність помилки). Чим вищий показник якості, тим надійнішим є показник розпізнавання. Q30 означає, що ймовірність помилкового розпізнавання в цій позиції становить 1 до 1000.
- 3.24 Read:** Одиниця безперервної послідовності ДНК (пар нуклеотидів), отримана шляхом секвенування частини фрагментованої ДНК-мішені.
- 3.25 Номер SAMN:** Унікальний ідентифікатор NCBI (номер доступу) для біозразка послідовності (метадані).
- 3.26 СОП:** Стандартна операційна процедура.
- 3.27 SRA: Архів** послідовного читання.
- 3.28 Номер SRR:** Унікальний ідентифікатор NCBI (номер приєднання) для завантажених зчитувань сирової послідовності.
- 3.29 Terra.Bio:** Хмарна платформа для аналізу послідовностей, розроблена Інститутом Броуда Массачусетського технологічного інституту та Гарвардським університетом і використовується компанією Theiagen Genomics (Highlands Ranch, Колорадо, США), щоб запропонувати платну спільну платформу для лабораторій громадського здоров'я для розміщення, аналізу та обміну даними (Libuit *et al.*, 2023).
- 3.30 TSV:** Значення, розділені табуляцією.
- 3.31 ПГС:** Повногеномне секвенування всього геному.

#### 4. **ОБОВ'ЯЗКИ:**

- 4.1 Персонал PulseNet виконує оцінку якості та генотипування коротких послідовностей зчитування Illumina, створених для спостереження PulseNet International з використанням платформи Terra.Bio. Наполегливо рекомендується ділитися даними WGS з іншими учасниками PNI шляхом завантаження в NCBI.

#### 5. **ПРОЦЕДУРА:**

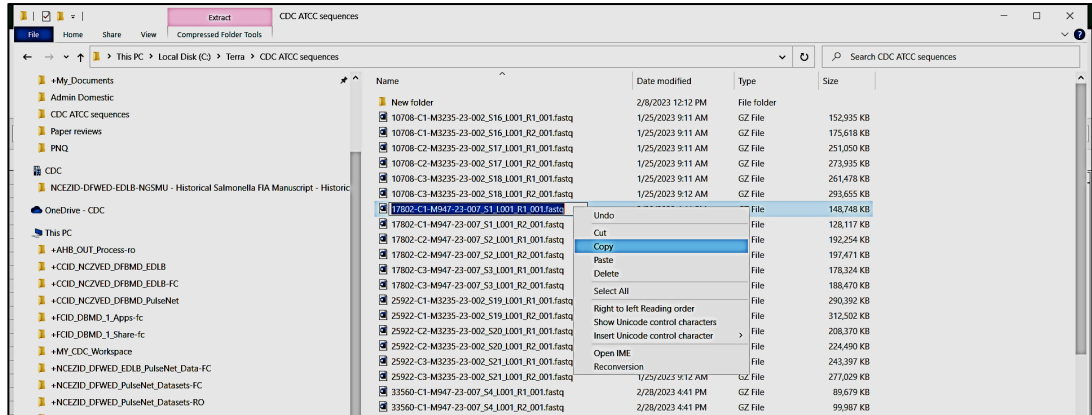
**ПРИМІТКА:** Дані послідовностей можна імпортувати в Terra одним з трьох способів: (1) завантаження з локального мережевого сховища (кроки 5.1 - 5.3), (2) пряма передача з хмари в хмару з Illumina BaseSpace ([додаток PNID01-1](#)), (3) завантаження з NCBI SRA ([додаток PNID01-2](#))

- 5.1 Підготуйте файл метаданих для послідовностей, які будуть завантажені на Terra . Файл метаданих пов'язуватиме завантажені FASTQ-файли з відповідними іменами записів у базі даних, тобто ключами записів.
- 5.1.1 Використовуйте Excel, щоб відкрити шаблон файлу метаданих у форматі tsv.

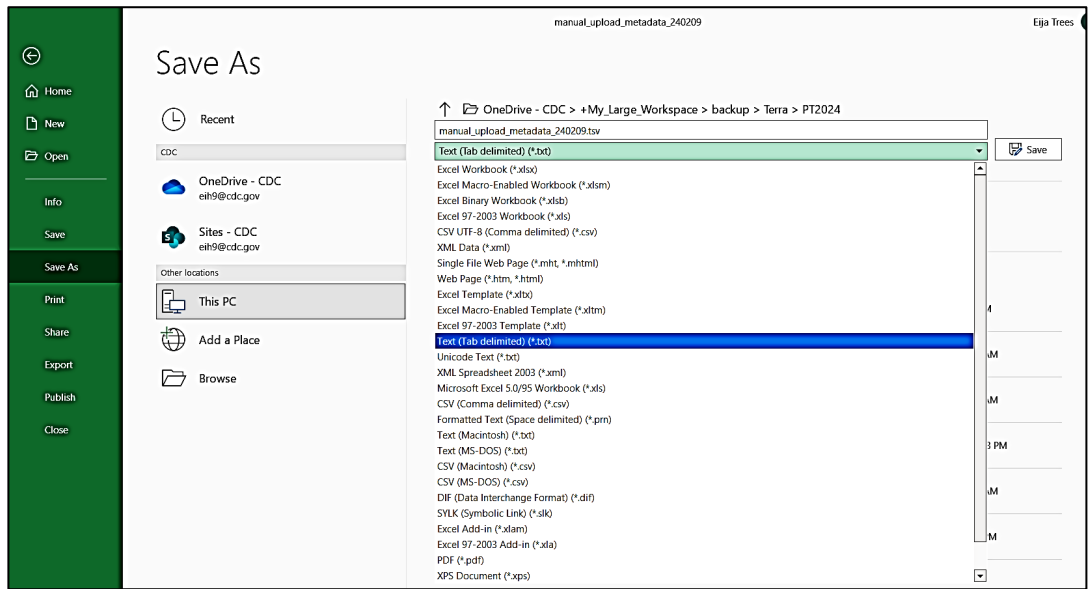


**ПРИМІТКА:** ідентифікатор штаму не обов'язково повинен збігатися з будь-якою частиною імені файлу *fastq.gz*.

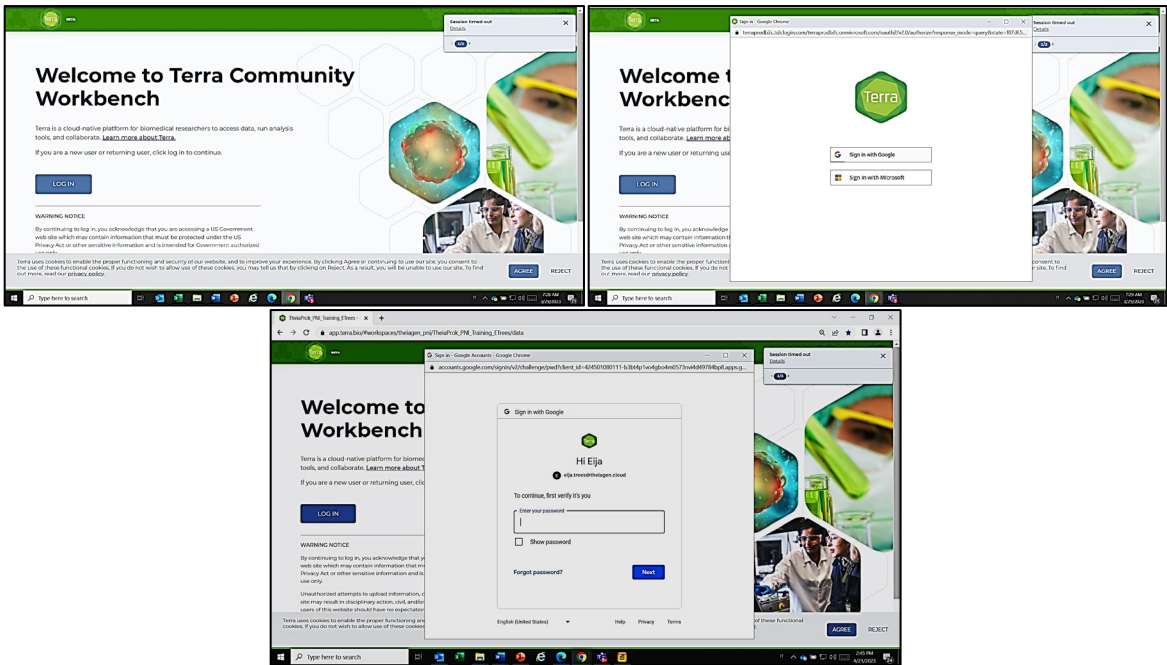
5.1.4 Скопіюйте та вставте імена файлів *fastq.gz* для кожного штаму в колонки "read1" та "read2". Переконайтеся, що імена файлів закінчуються на "fastq.gz".



5.1.5 Збережіть файл у форматі **tsv**: виберіть "Зберегти як" і "Текст (з розділенням табуляцією) (\*.txt)". Переконайтеся, що ім'я вашого файлу закінчується на ".tsv".

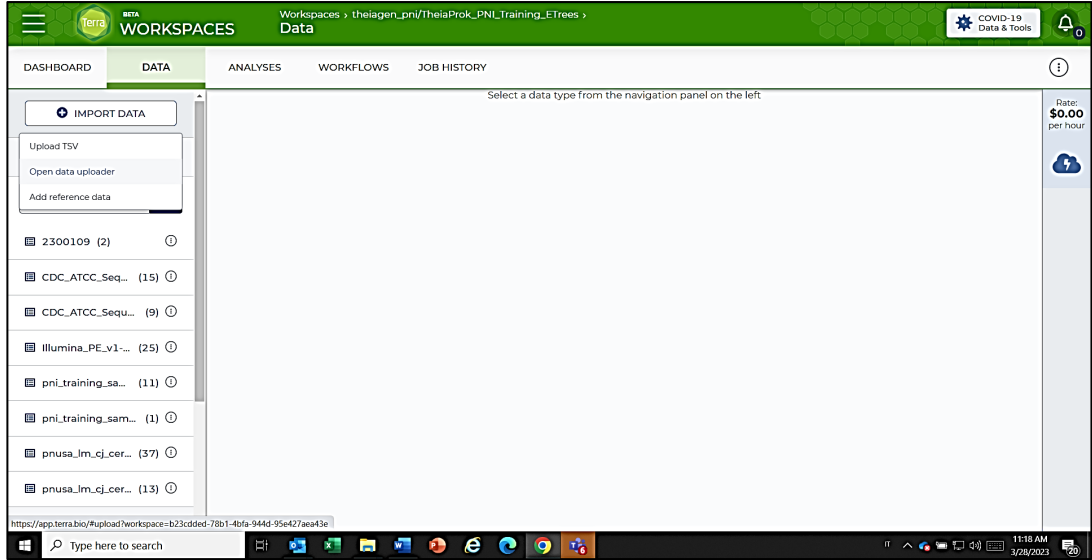


5.2 Увійдіть на Terra за допомогою Chrome та свого облікового запису Google

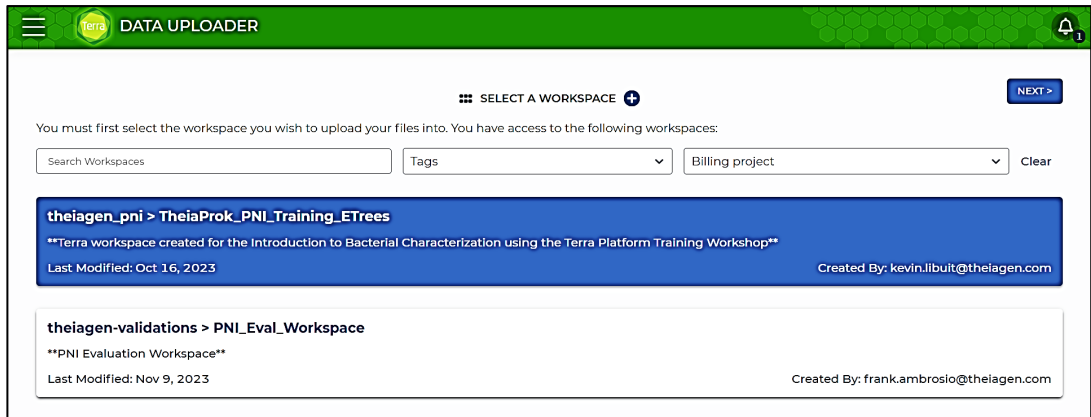


**5.3 Завантажте файли послідовностей та метадані до Terra**

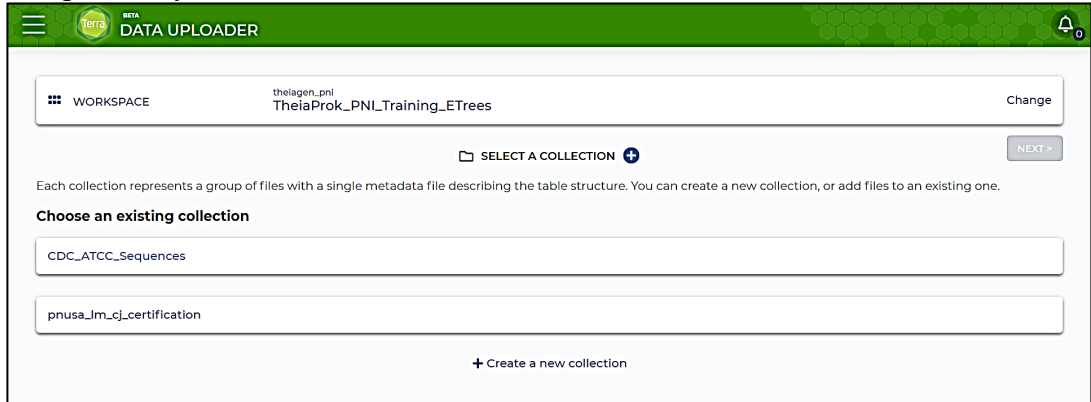
5.3.1 У розділі "Робочі простори Terra" виберіть вкладку "Дані", натисніть "Імпортувати дані" і виберіть "Відкрити завантажувач даних" у випадаючому МЕНЮ.



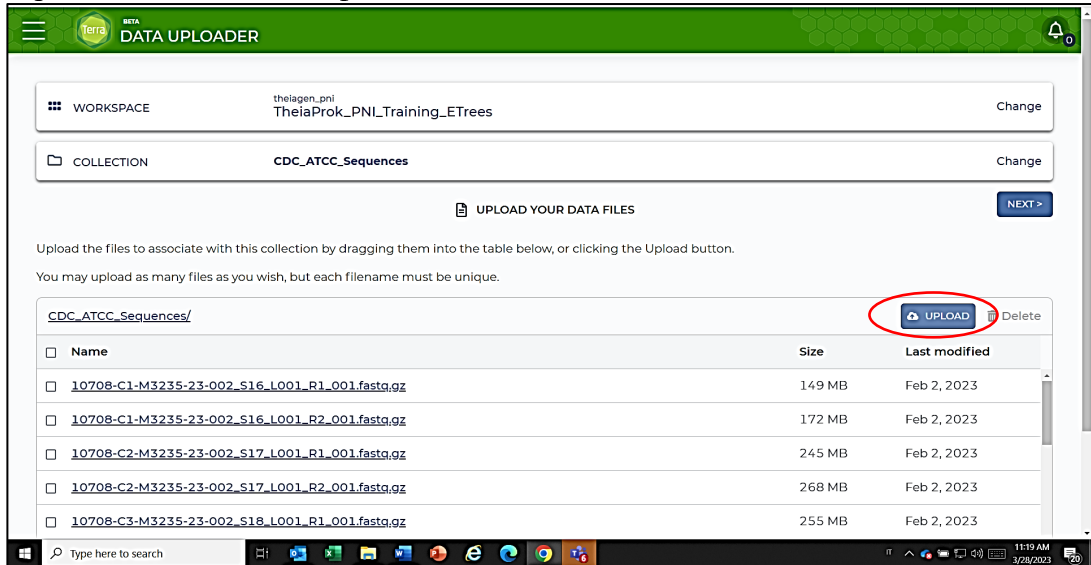
5.3.2 Відкриється вікно "Завантажувач даних". Якщо ваш обліковий запис має доступ до кількох робочих просторів, спочатку потрібно вибрати робочий простір, до якого ви бажаєте завантажити дані.



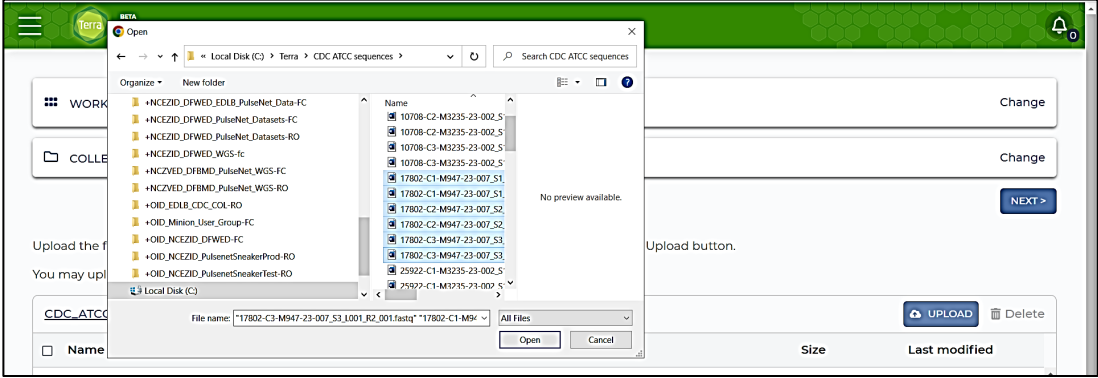
5.3.3 Або виберіть існуючу колекцію, натиснувши на назву колекції у списку, або створіть нову колекцію.



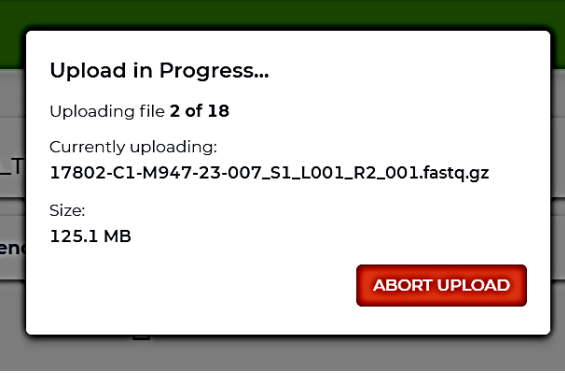
5.3.4 У розділі "Завантажити файли даних" натисніть "Завантажити".



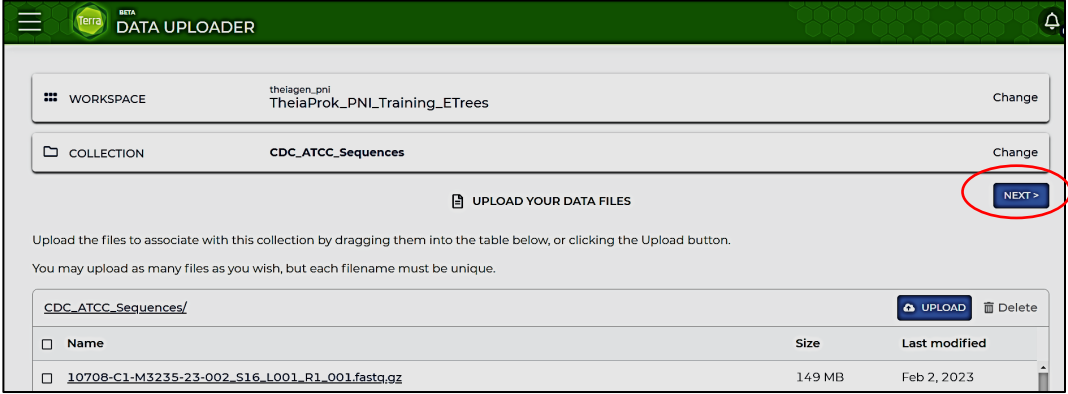
5.3.5 Перейдіть до місця, де зберігаються файли FASTQ, виберіть файли, які потрібно завантажити, і натисніть "Відкрити".



5.3.6 З'явиться спливаюче вікно "Завантаження триває". Завантаження може зайняти кілька хвилин, залежно від кількості файлів, що завантажуються, та пропускну здатності вашого Інтернету.

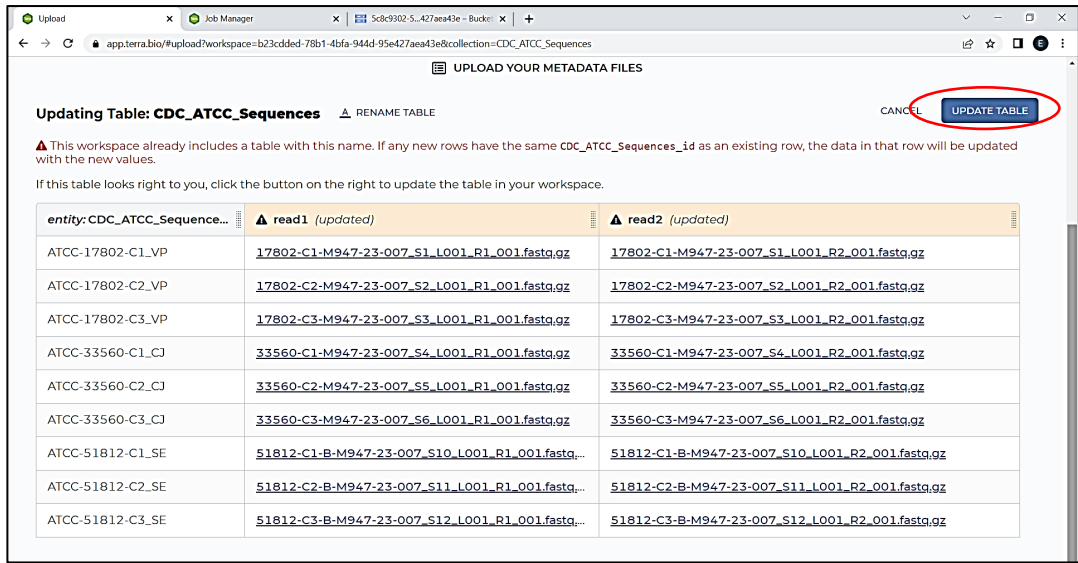


5.3.7 Після того, як спливаюче вікно зникне, натисніть "Далі".



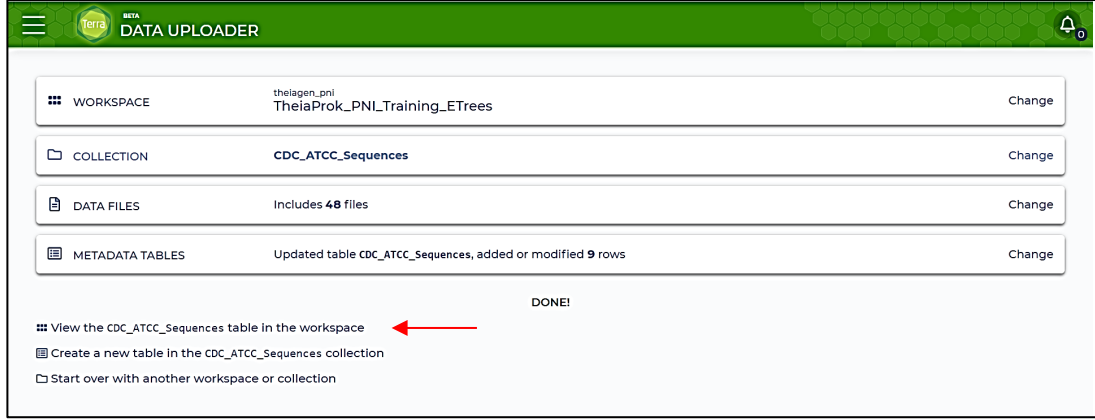
5.3.8 У розділі "Завантажити файли метаданих" натисніть "Завантажити"





5.3.11 На екрані "Завантаження даних" з'явиться повідомлення "Готово". Ви можете переглянути оновлену таблицю даних, натиснувши на посилання, яке з'явиться на екрані.

**ПРИМІТКА:** *Стовпці, що відображаються в таблиці даних, можна налаштувати. Для зручності навігації рекомендується створити окремі подання для показників КК, результатів генотипування та метаданих. Зверніться до додатків [PNID01-3](#) (Показники контролю якості), [PNID01-5](#) (Генотипування) та [PNID01-6](#) (Метадані) для отримання інструкцій щодо налаштування стовпців таблиці даних для спостереження PulseNet.*



**МІЖНАРОДНА СТАНДАРТНА ОПЕРАЦІЙНА ПРОЦЕДУРА PULSENET ДЛЯ АНАЛІЗУ ДАНИХ  
КОРОТКОГО ЗЧИТУВАННЯ WGS ILLUMINA З ВИКОРИСТАННЯМ ПЛАТФОРМИ TERRA.BI**

Док. No. PNID01

Вер. No. 01

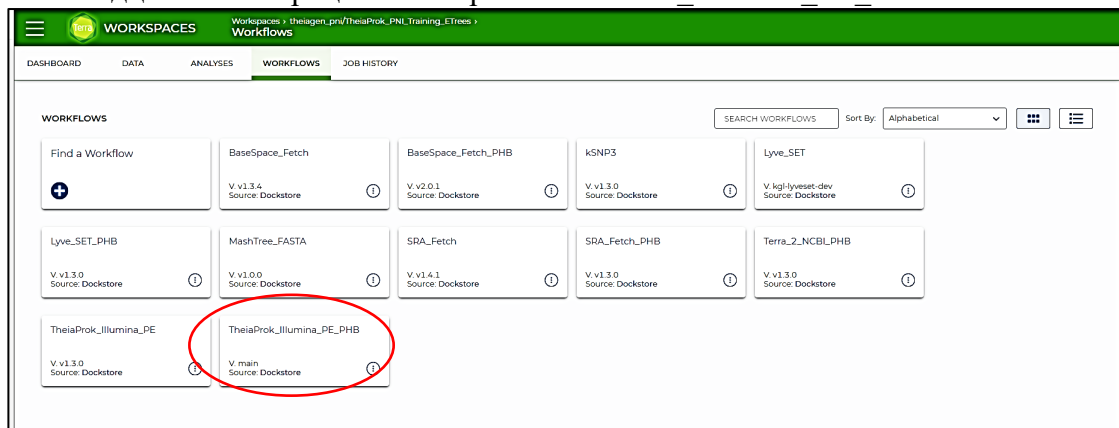
Дата набуття чинності:

Сторінка 11 з 67

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

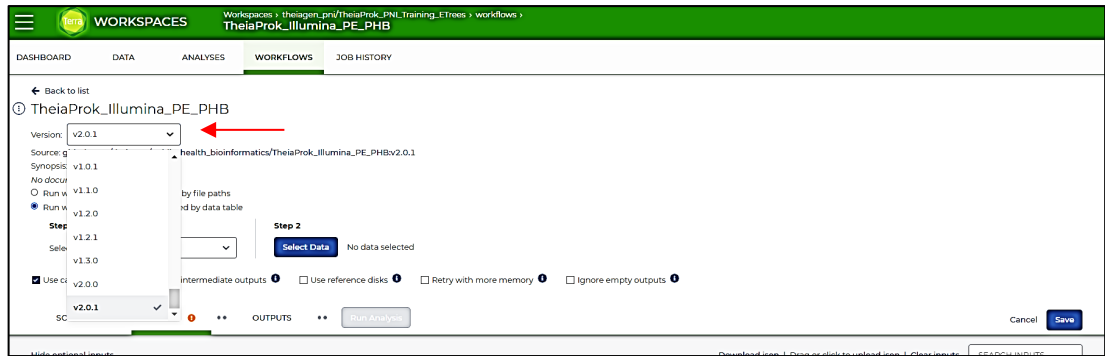
5.4 **Запустіть робочий процес контролю якості та генотипування.** Робочий процес TheiaProk виконує QC сирих зчитувань послідовностей, збирає сирі зчитування de novo за допомогою Skesa, а потім виконує QC збірки та ідентифікацію виду. Також доступні різноманітні аналізи генотипування, що підходять для даного виду.

5.4.1 На вкладці "Робочі процеси" виберіть "TheiaProk\_Illumina\_PE\_PHB".



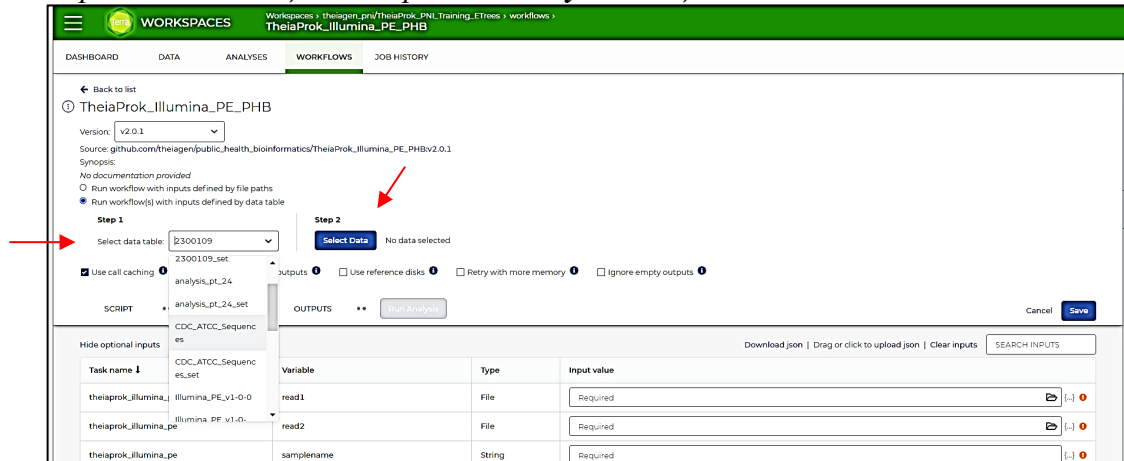
5.4.2 На екрані "TheiaProk\_Illumina\_PE\_PHB":

5.4.2.1 Виберіть останню версію робочого процесу TheiaProk\_Illumina\_PE\_PHB зі спадного меню "Версія".



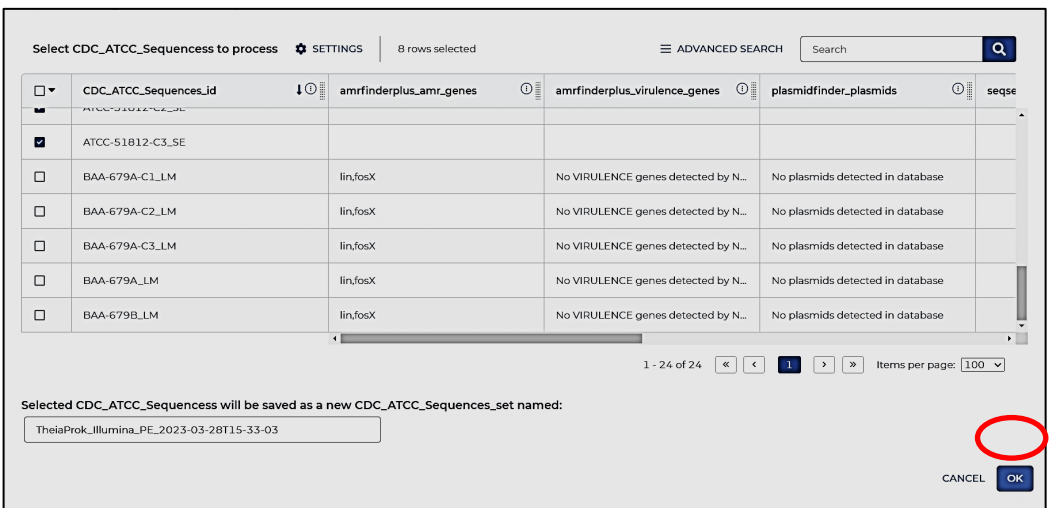
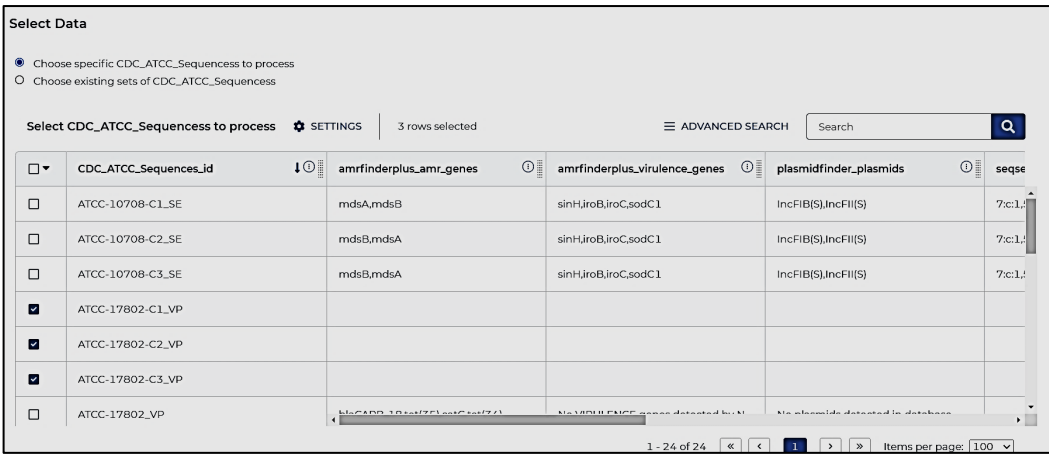
5.4.2.2 У розділі "Крок 1" виберіть "Тип кореня", тобто таблицю даних, в якій знаходяться зразки, наприклад, "CDC\_ATCC\_Sequences".

**ПРИМІТКА:** Для більшості таблиць даних є два варіанти: *основна таблиця даних, що містить окремі записи зразків, і таблиця даних наборів, що містить набори зразків, які використовуються для філогенетичного аналізу, завантаження з NCBI тощо (наприклад, CDC\_ATCC\_Sequences\_set).* *Переконайтеся, що ви вибрали головну таблицю даних.*

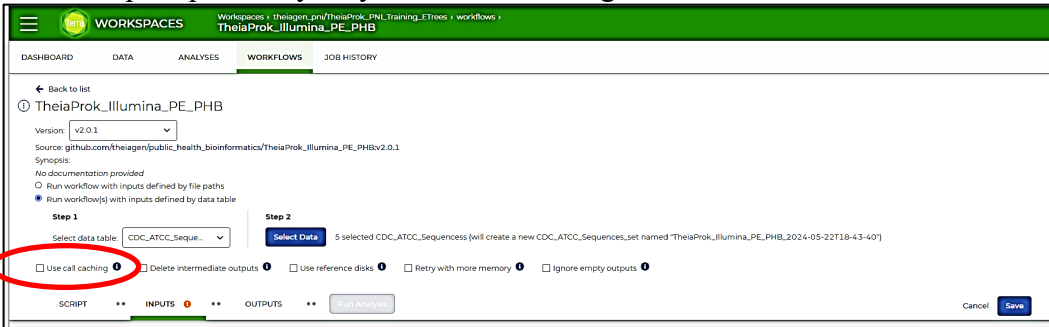


5.4.2.3 У розділі "Крок 2" натисніть "Вибрати дані" (скріншот вище). Ви потрапите до таблиці даних, вказаної на кроці 1.

5.4.2.4 Виберіть штами для аналізу і прокрутіть сторінку донизу, щоб натиснути "ОК". **ПРИМІТКА:** якщо таблиця даних містить понад 100 записів і ви встановите прапорець "Виділити все" біля назви таблиці даних, буде вибрано лише 100 записів.



5.4.2.5 Зніміть прапорець з пункту "Use call catching".



5.4.2.6 На вкладці "Вхідні дані" вкажіть наступні вхідні значення в колонці "Атрибут" (прокрутіть список вниз):

**ПРИМІТКА:** Коли ви заповнюєте стовпчик "Атрибут", клацання всередині клітинки відкриває спадне меню атрибутів, які ви можете вибрати, щоб уникнути помилок (скріншот нижче).

5.4.2.6.1 Зчитування 1: "This.read1".

5.4.2.6.2 Read2: "This.read2".

5.4.2.6.3 Ім'я зразка: "**This.data table name\_id**", наприклад, This.CDC\_ATCC\_Sequences\_id.

5.4.2.6.4 Call\_ani: "True".

Task name ↓	Variable	Type	Input value
theiaprok_illumina_pe	read1	File	this.read1
theiaprok_illumina_pe	read2	File	this.read2
theiaprok_illumina_pe	samplename	String	this.CDC_
amrfinderplus_task	cpu	int	this.CDC_ATCC_Sequences_id
amrfinderplus_task	detailed_drug_class	Boolean	Optional
shovill_pe	trim	Boolean	Optional
theiaprok_illumina_pe	call_ani	Boolean	true
theiaprok_illumina_pe	call_kmerfinder	Boolean	Optional

5.4.2.7 На вкладці "**Outputs**" натисніть "Використовувати за замовчуванням".

Task name ↓	Variable	Type	Input value   Use defaults
theiaprok_illumina_pe	abricate_abaum_plasmid_tsv	File	this.abricate_abaum_plasmid_tsv
theiaprok_illumina_pe	abricate_abaum_plasmid_type_genes	String	this.abricate_abaum_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 Натисніть "Зберегти" (скріншот вище).

**ПРИМІТКА:** кнопка "Зберегти" відображається лише в тому випадку, якщо ви змінили вхідні дані з попереднього подання.

5.4.2.9 Натисніть "Запустити аналіз". З'явиться спливаюче вікно "Підтвердити запуск", в якому ви можете ввести необов'язковий опис. Натисніть "Запустити".

Run workflow(s) with inputs defined by the table

**Step 1**  
 Select data table: CDC\_ATCC\_Sequences...  
 Use call caching  Delete intermediate outputs  Use reference disks  Retry with more memory  Ignore empty outputs

**Step 2**  
 Select Data 3 selected CDC\_ATCC\_Sequences (will create a new CDC\_ATCC\_Sequences\_set named ...)

Output files will be saved to: Files / submission unique ID / theiaprok\_illumina\_pe / workflow unique ID  
 References to outputs will be written to: Tables / CDC\_ATCC\_Sequences

Task name ↓	Variable	Type	Input value   Use defaults
theiaprok_illumina_pe	abricate_abaum_plasmid_tsv	File	this.abricate_abaum_plasmid_tsv

**Confirm launch**

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

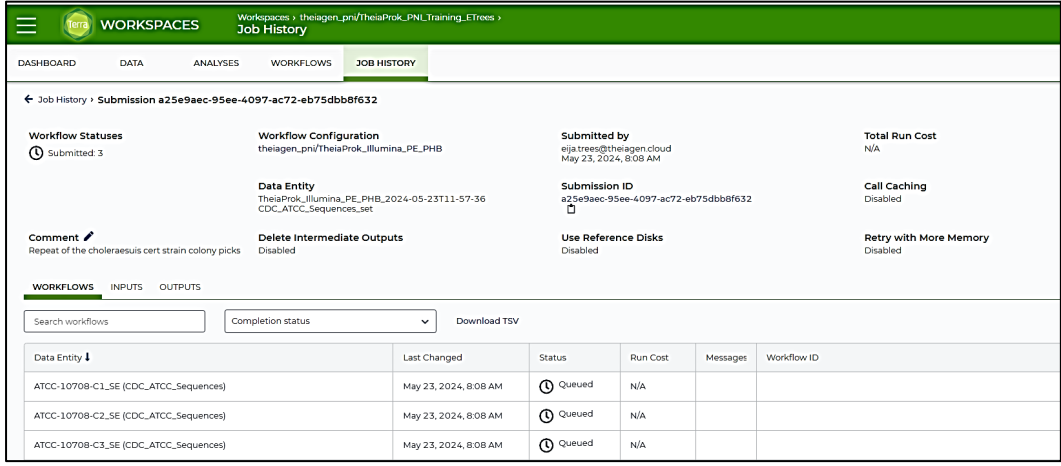
Running workflows will generate cloud charges. How much does my workflow cost? Set up budget alert

Describe your submission (optional):  
 Colony picks for VP, C3 and ST1

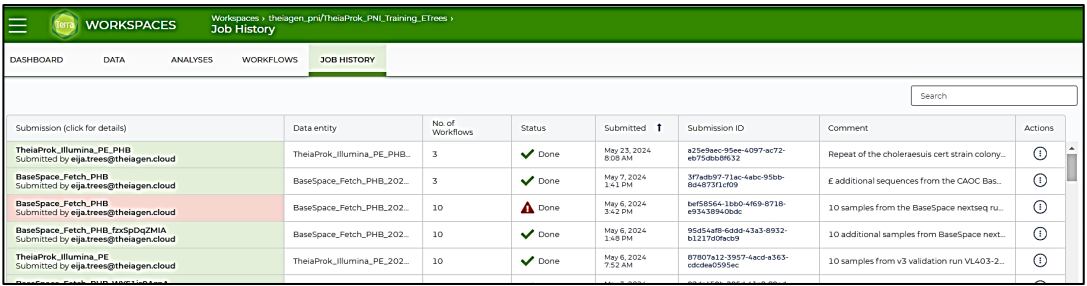
This will launch 8 analyses.

CANCEL LAUNCH

5.4.3 З'явиться вікно "Статуси робочого процесу", де ваші надіслані завдання повинні бути спочатку перераховані як "В черзі".

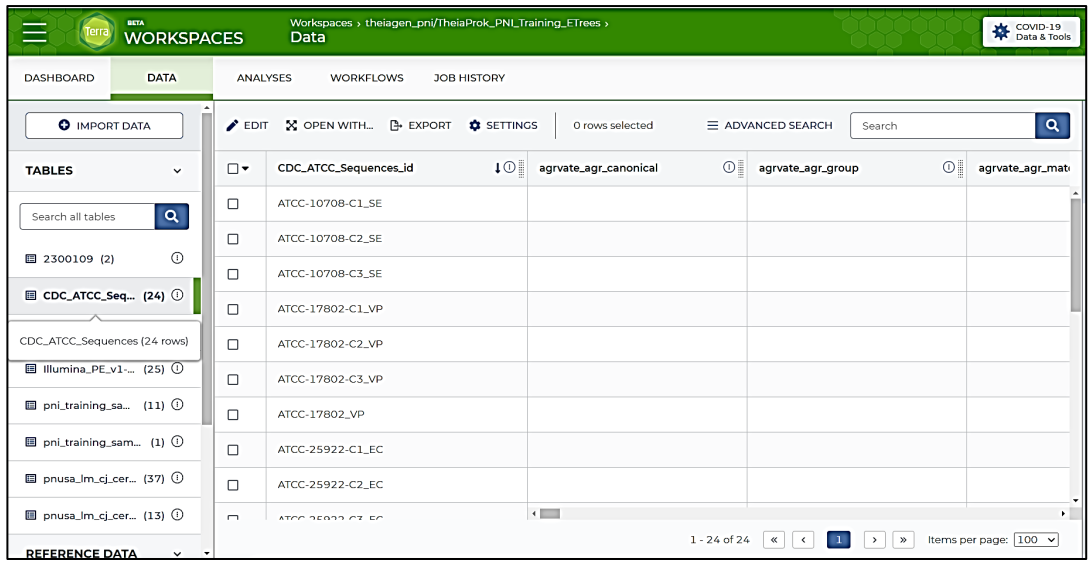


5.4.4 Перейдіть на вкладку "Історія завдань", щоб перевірити статус вашого завдання. Успішно виконане завдання позначається зеленою галочкою.

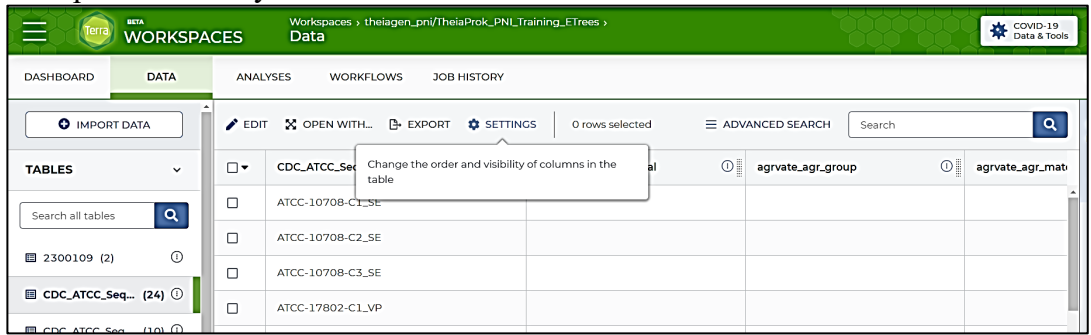


5.5 Оцініть метрики КЯ для послідовностей: Показники КЯ можна переглянути безпосередньо в таблиці даних (5.5.1-5.5.3) або експортувати в Excel для вибраних записів (5.5.4).

5.5.1 У "Робочих просторах Terra" виберіть вкладку "Дані", а потім виберіть таблицю даних, яка вас цікавить, наприклад, "CDC\_ATCC\_Sequences".



5.5.2 Виберіть "Налаштування".



5.5.3 На екрані "Вибір стовпців" у розділі "Ваші збережені вибори стовпців" натисніть на коло з 3 крапками поруч з "qc\_metrics" і у випадяючому меню виберіть "Завантажити", а потім натисніть "Готово". Це призведе до завантаження відповідних метрик PulseNet QC в таблицю даних. Зверніться до Додатку [PNID01-3](#) для отримання інформації про показники QC, які повинні з'явитися в таблиці, а також для отримання інструкцій про те, як додати або видалити будь-який з стовпців (показники QC) в таблиці показників QC.

**МІЖНАРОДНА СТАНДАРТНА ОПЕРАЦІЙНА ПРОЦЕДУРА PULSENET ДЛЯ АНАЛІЗУ ДАНИХ  
КОРОТКОГО ЗЧИТУВАННЯ WGS ILLUMINA З ВИКОРИСТАННЯМ ПЛАТФОРМИ TERRA.VI**

Док. No. PNID01

Вер. No. 01

Дата набуття чинності:

Сторінка 17 з 67

**Select columns**

Show: all | none      Sort: alphabetical

- amrfinderplus\_all\_report
- 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

—

SAVE THIS COLUMN SELECTION

Your saved column selections:

- pulsenet\_genotyping ⓘ
- qc\_metrics ⓘ ←

qc\_metrics

CANCEL    **DONE**

**Select columns**

Show: all | none      Sort: alphabetical

- agrvate\_agr\_canonical
- agrvate\_agr\_group
- agrvate\_agr\_match\_score
- agrvate\_agr\_multiple
- agrvate\_agr\_num\_frameshifts
- agrvate\_docker
- agrvate\_results
- agrvate\_summary
- agrvate\_version
- meningotype\_BAST
- meningotype\_FetA
- meningotype\_NHBA
- meningotype\_NadA

—

SAVE THIS COLUMN SELECTION

Your saved column selections:

- pulsenet\_genotyping ⓘ
- qc\_metrics ⓘ

Load

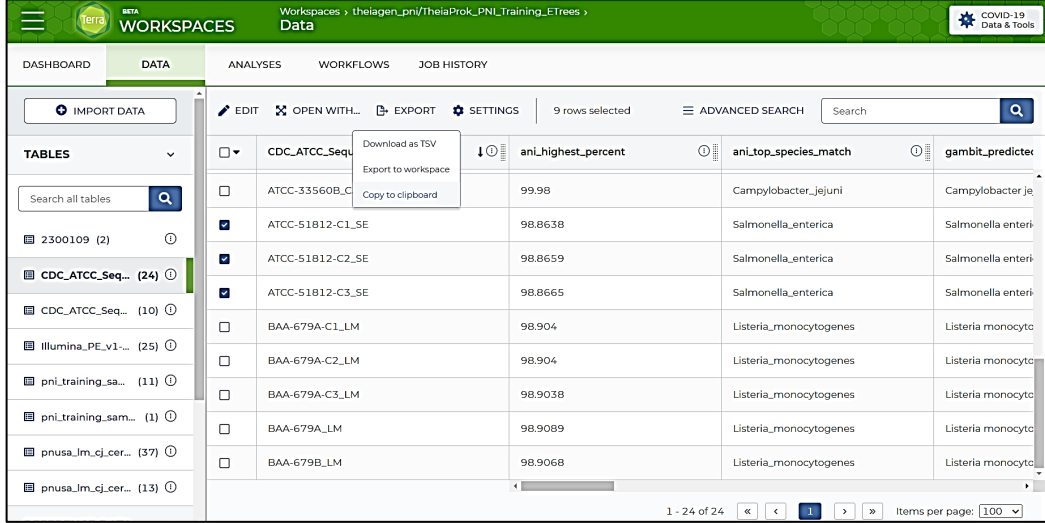
Delete

CANCEL    **DONE**

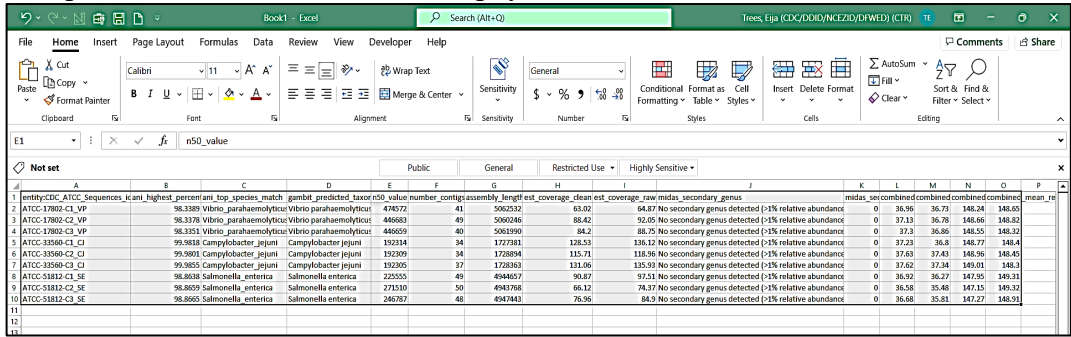
5.5.4 Експорт метрик контролю якості в Excel:

5.5.4.1 Виберіть потрібні записи.

5.5.4.2 Натисніть "Експортувати", а потім виберіть "Копіювати в буфер обміну".



5.5.4.3 Відкрийте Excel і вставте дані на аркуш.



5.5.5 Критичні показники якості PulseNet для прийнятних послідовностей Illumina див. у Додатку [PNID01-4a](#).

5.5.6 TheiaProk використовує етап попередньої перевірки зчитування, на якому він зупиняє аналіз для зразків, які не досягають певних порогових значень якості. У цьому випадку на вкладці "Історія завдань" буде вказано, що зразок пройшов успішно, але на вкладці "Дані" не буде результатів. У колонці "raw\_read\_screen" в матриці qc має бути вказана причина помилки аналізу послідовності. Порогові значення, що застосовуються на цьому етапі попереднього скринінгу зчитування, наведено в Додатку [PNID01-4b](#).

**МІЖНАРОДНА СТАНДАРТНА ОПЕРАЦІЙНА ПРОЦЕДУРА PULSENET ДЛЯ АНАЛІЗУ ДАНИХ КОРОТКОГО ЗЧИТУВАННЯ WGS ILLUMINA З ВИКОРИСТАННЯМ ПЛАТФОРМИ TERRA.VI**

Док. No. PNID01

Вер. No. 01

Дата набуття чинності:

Сторінка 19 з 67

The screenshot shows the Terra Workspaces interface. The top navigation bar includes 'WORKSPACES' and 'Job History'. The main content area displays details for a workflow submission with ID 'c4ba38f7-33b1-4e05-ac0c-8aa3c742f8b3'. Key information includes:
 

- Workflow Statuses:** Succeeded (1)
- Workflow Configuration:** theiaigen\_pnl/TheiaProk\_PNL\_Training\_ETrees
- Submitted by:** eja.trees@theiaigen.cloud, Mar 25, 2024, 9:04 AM
- Total Run Cost:** N/A
- Data Entity:** 2013L-5615TK\_NextSeq\_400MB analysis\_pt\_24
- Submission ID:** c4ba38f7-33b1-4e05-ac0c-8aa3c742f8b3
- Call Caching:** Enabled
- Comment:** Repeat of 2013L-5615
- Delete Intermediate Outputs:** Disabled
- Use Reference Disks:** Disabled
- Retry with More Memory:** Disabled

 Below the details is a table with columns: Data Entity, Last Changed, Status, Run Cost, Message, and Workflow ID. One row is visible for '2013L-5615TK\_NextSeq\_400MB (analysis\_pt\_24)' with a status of 'Succeeded'.

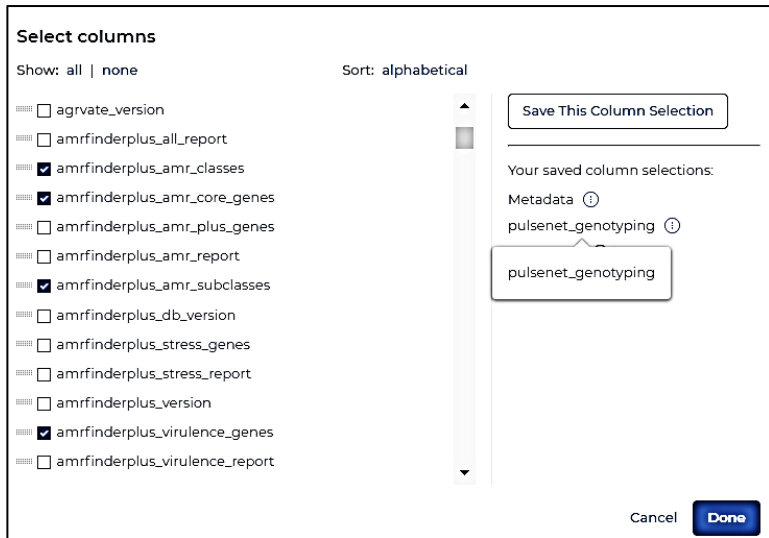
The screenshot shows the Terra Workspaces 'Data' table. The table has columns: analysis\_pt\_24\_id, midas\_secondary.genus, midas\_secondary.genus.abundanc, n50.value, number\_contigs, and raw\_read\_screen. The table contains 10 rows of data. The 9th row is highlighted in blue and shows a 'FAIL' status due to low coverage.

analysis_pt_24_id	midas_secondary.genus	midas_secondary.genus.abundanc	n50.value	number_contigs	raw_read_screen
2011V-1043_FLEX_300_Vbribo	No secondary genus detected (>1% r...	0	124949	77	PASS
2012V-1116_FLEX_300_Vbribo	No secondary genus detected (>1% r...	0	458045	39	PASS
2013L-5361_FLEX_300_LM	No secondary genus detected (>1% r...	0.0006	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	No secondary genus detected (>1% r...	0	443204	15	FAIL: the estimated coverage is less than the minimum of 10x
2015AM-1304	No secondary genus detected (>1% r...	0	728098	16	PASS
2015AM-1305	No secondary genus detected (>1% r...	0	728098	16	PASS

**5.6 Перегляньте результати генотипування для послідовностей: Результати**

генотипування можна переглянути або безпосередньо в таблиці даних (5.6.1-5.6.3), або експортувати в Excel для вибраних записів (5.6.4).

- 5.6.1 У "Робочих просторах Terra" виберіть вкладку "Дані", а потім виберіть таблицю даних, яка вас цікавить, наприклад, "CDC\_ATCC\_Sequences".
- 5.6.2 На вкладці "Дані" виберіть "Налаштування".
- 5.6.3 На екрані "Вибір стовпчиків" у розділі "Ваш збережений вибір стовпчиків" натисніть на коло з 3 крапками поруч з "pulsenet\_genotyping" і у випадяючому меню виберіть "Завантажити", а потім натисніть "Готово". Це призведе до завантаження аналізів генотипування, які підходять для спостереження PulseNet, в таблицю даних. Зверніться до Додатку [PNID01-5](#) для отримання інформації про тести генотипування, які повинні з'явитися в таблиці, а також для отримання інструкцій, як додати або видалити будь-яку колонку (тести генотипування) в таблиці результатів генотипування.



5.6.4 Експортуйте результати в Excel для вибраних записів; виконайте процедуру, описану в кроці 5.5.4.

**5.7 Завантажуйте послідовності до NCBI**

**ПРИМІТКА:** Зверніться до компанії *Theiagen Genomics* ([support@theiagen.com](mailto:support@theiagen.com)) для отримання інструкцій перед початком роботи та налаштування робочого простору для завантаження даних до NCBI. Процес конфігурації описано за посиланням: [https://theiagen.notion.site/Terra\\_2\\_NCBI-61abcedc066646b3b258f70b561e9f62](https://theiagen.notion.site/Terra_2_NCBI-61abcedc066646b3b258f70b561e9f62).

- 5.7.1 Завантажте метадані, необхідні для подання до NCBI: див. [додаток PNID01-6](#) для правильного форматування та завантаження метаданих.
- 5.7.2 Створіть **набір** зразків для завантаження в NCBI:
  - 5.7.2.1 У розділі "Робочі простори Terra" виберіть вкладку "Дані", а потім виберіть таблицю даних, яка вас цікавить, наприклад, "контроль якості".
  - 5.7.2.2 На вкладці "Дані" виберіть послідовності, які потрібно включити до завантаження в NCBI.
  - 5.7.2.3 У випадваючому меню "Редагувати" виберіть "Зберегти виділення як задане".

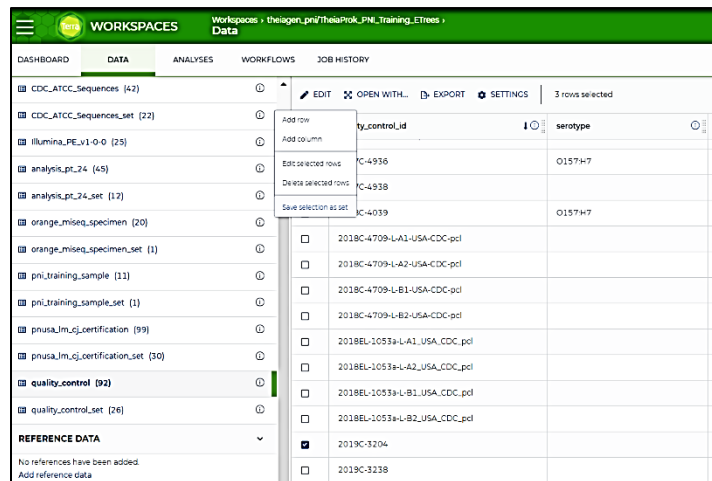
**МІЖНАРОДНА СТАНДАРТНА ОПЕРАЦІЙНА ПРОЦЕДУРА PULSENET ДЛЯ АНАЛІЗУ ДАНИХ КОРОТКОГО ЗЧИТУВАННЯ WGS ILLUMINA З ВИКОРИСТАННЯМ ПЛАТФОРМИ TERRA.VI**

Док. No. PNID01

Вер. No. 01

Дата набуття чинності:

Сторінка 21 з 67



5.7.2.4 У спливаючому вікні назвіть набір, наприклад, "NCBI\_upload\_240422" і натисніть "Зберегти".

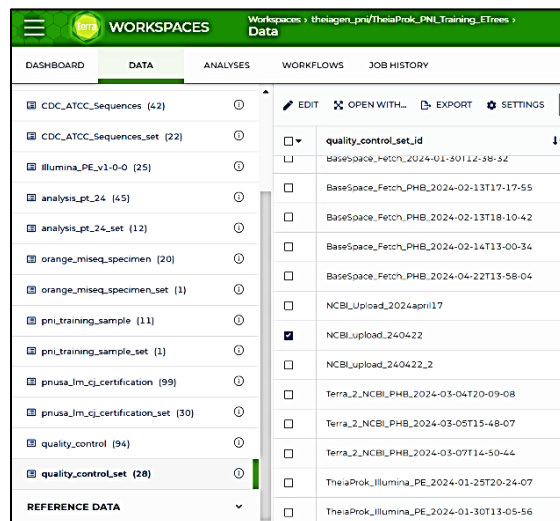
**ПРИМІТКА:** пробіли та тире не допускаються.

Create a quality\_control set

Set name (required)

CANCEL **SAVE**

5.7.2.5 Новостворений набір тепер має з'явитися в таблиці даних "набір", наприклад, "набір\_контролю\_якості".



5.7.3 Налаштуйте параметри робочого процесу завантаження до NCBI:

5.7.3.1 На вкладці "Робочі процеси" виберіть робочий процес "Terra\_2\_NCBI\_PHB". Відкриється вікно "Terra\_2\_NCBI\_PHB".

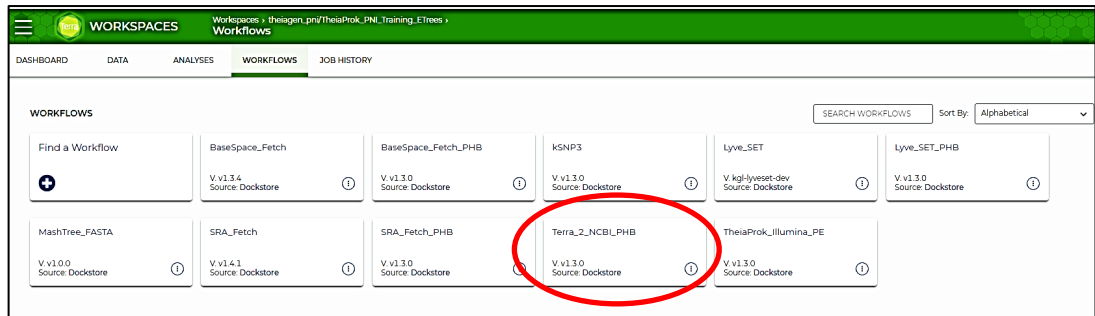
**МІЖНАРОДНА СТАНДАРТНА ОПЕРАЦІЙНА ПРОЦЕДУРА PULSENET ДЛЯ АНАЛІЗУ ДАНИХ  
КОРОТКОГО ЗЧИТУВАННЯ WGS ILLUMINA З ВИКОРИСТАННЯМ ПЛАТФОРМИ TERRA.VI**

Док. No. PNID01

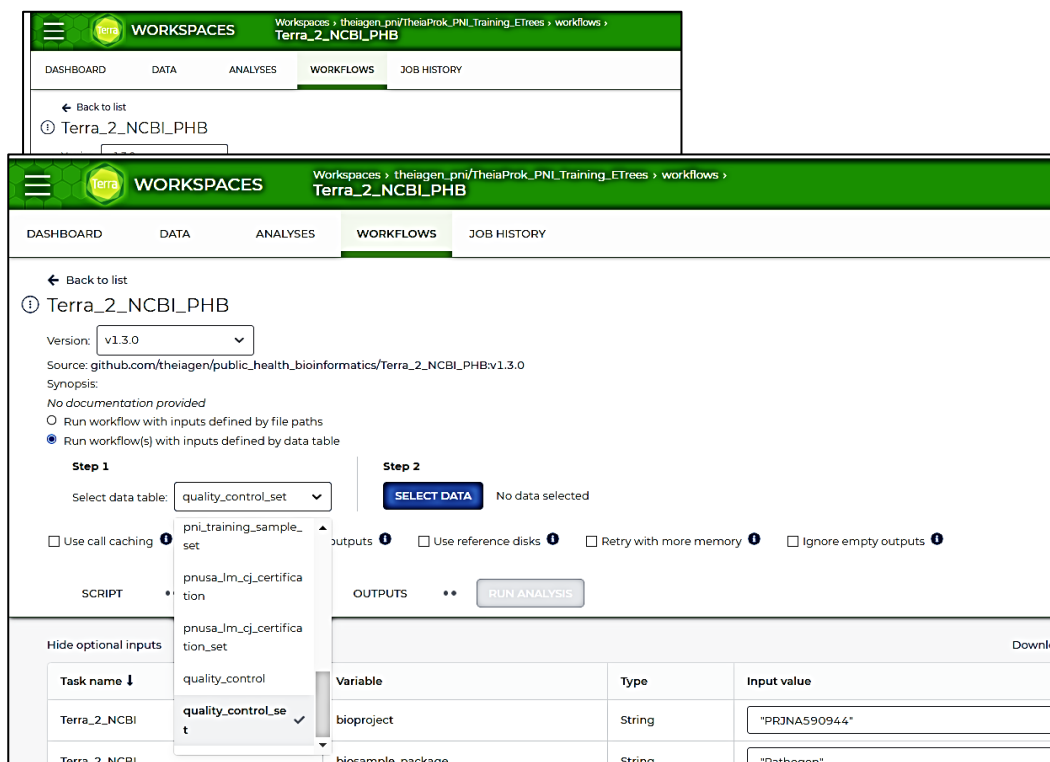
Вер. No. 01

Дата набуття чинності:

Сторінка 22 з 67



5.7.3.2 У випадяючому меню "Версія" виберіть останню версію "Terra\_2\_NCBI\_PHB".



5.7.3.3 На кроці 1 у випадяючому меню "Select root entity type" виберіть таблицю даних набору, в якій знаходиться зразок набору, створений на кроці 5.7.2, наприклад, "quality\_control\_set". Також зніміть позначку з пункту " Use call caching ".

5.7.3.4 На кроці 2 натисніть "Вибрати дані" (скріншот вище). Ви потрапите до таблиці даних, вибраної на попередньому кроці.

5.7.3.5 Виберіть потрібний набір зразків, наприклад, "NCBI\_upload\_240422" і натисніть "OK".

**МІЖНАРОДНА СТАНДАРТНА ОПЕРАЦІЙНА ПРОЦЕДУРА PULSENET ДЛЯ АНАЛІЗУ ДАНИХ  
КОРОТКОГО ЗЧИТУВАННЯ WGS ILLUMINA З ВИКОРИСТАННЯМ ПЛАТФОРМИ TERRA.VI**

Док. No. PNID01

Вер. No. 01

Дата набуття чинності:

Сторінка 23 з 67

Select Data

Create a new quality\_control\_set from selected quality\_controls  
 Choose specific quality\_control\_sets to process

Select quality\_control\_sets to process **SETTINGS** | 1 row selected

ADVANCED SEARCH Search

<input type="checkbox"/>	quality_control_set_id	Terra_2_NCBL_analysis_date	Terra_2_NCBL_version	biosample_failures	biosample_metadata	biosample_report_xmls
<input type="checkbox"/>	BaseSpace_Fetch_PHB_2024-04-22T13:58:04					
<input type="checkbox"/>	NCBL_upload_2024apr117	2024-04-17	PHB v1.3.0	biosample_failures.txt	biosample_table.tsv	biosample_table-report.2
<input checked="" type="checkbox"/>	NCBL_upload_240422	2024-04-22	PHB v1.3.0	biosample_failures.txt	biosample_table.tsv	biosample_table-report.2
<input type="checkbox"/>	NCBL_upload_240422_2	2024-04-24	PHB v1.3.0		biosample_table.tsv	
<input type="checkbox"/>	Terra_2_NCBL_PHB_2024-03-04T20-09-08	2024-03-04	PHB v1.3.0	biosample_failures.txt	biosample_table.tsv	biosample_table-report.2
<input type="checkbox"/>	Terra_2_NCBL_PHB_2024-03-05T15-48-07	2024-03-05	PHB v1.3.0	biosample_failures.txt	biosample_table.tsv	biosample_table-report.2
<input type="checkbox"/>	Terra_2_NCBL_PHB_2024-03-07T14-50-44	2024-03-07	PHB v1.3.0	biosample_failures.txt	biosample_table.tsv	biosample_table-report.2

1 - 28 of 28    Items per page:

CANCEL OK

5.7.3.6 На вкладці "Вхідні дані" потрібно заповнити наступні "Вхідні значення" для перелічених нижче "Змінних":

5.7.3.6.1 біопроект: введіть номер біопроекту NCBI, до якого ви бажаєте завантажити дані, у лапках, наприклад, "PRJNA590944".

5.7.3.6.2 biosample\_package в лапках: "Патоген". Це шаблон/пакет метаданих, який ви використовуєте для завантаження метаданих для послідовностей, що належать до спостереження PulseNet.

5.7.3.6.3 ncbi\_config\_js: введіть назву файлу конфігурації NCBI, створеного для вашої робочої області, наприклад, workspace.ncbi\_config\_etrees.

5.7.3.6.4 назва\_проекту в лапках: "theiagen\_pni".

5.7.3.6.5 sample\_names: введіть назву вашої таблиці даних у форматі:

**this.data\_table\_names.data\_table\_name\_id,**

наприклад, це\_якість\_контролює.ідентифікатор\_контролю\_якості.

**ПРИМІТКА:** формат подвійного імені є **ОБОВ'ЯЗКОВИМ**.

5.7.3.6.6 sra\_transfer\_gcp\_bucket в лапках: "gs://theiagen\_sra\_transfer". Це тимчасове загальнодоступне сховище Google для ваших послідовностей, до якого може отримати доступ NCBI.

5.7.3.6.7 назва\_таблиці: введіть назву вашої таблиці даних у лапках, наприклад, "quality\_control".

5.7.3.6.8 ім'я\_робочого\_простору: введіть ім'я вашого робочого простору в лапках, наприклад, "TheiaProk\_PNI\_Training\_ETrees".

5.7.3.6.9 submit\_to\_production: true.

5.7.3.6.10 НЕОБОВ'ЯЗКОВО: якщо ви хочете пов'язати заявку SRA з **існуючим біозразком** (поле "biosample\_accession" вже заповнене номером SAMN на вкладці "Дані");

5.7.3.6.10.1 skip\_biosample: true.

**ПРИМІТКА:** для використання цієї функції необхідно заповнити необхідні поля метаданих, тобто заповнення лише поля "biosample\_accession" є недостатнім і не пройде перевірку попереднього завантаження в Terra.

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 На вкладці "Вихідні дані" натисніть "Використовувати значення за замовчуванням" для "Атрибутів", а потім натисніть "Зберегти". **ПРИМІТКА:** кнопка "Зберегти" відображається, тільки якщо ви змінили вхідні дані з попереднього подання.

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_xmIs	Array[File]	this.biosample_report_xmIs
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 Натисніть "Запустити аналіз" (скріншот вище). З'явиться спливаюче вікно "Підтвердити запуск", в якому ви можете ввести необов'язковий опис. Натисніть "Запустити".

**Confirm launch**

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

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

5.7.3.9 З'явиться вікно "Статуси робочого процесу", де ваші надіслані завдання повинні бути спочатку перераховані як "В черзі" або "Запуск".

The screenshot shows the 'Job History' page for a submission with ID 387be0c7-88b6-4c00-8028-9020dc5f08b8. The page is divided into several sections:

- Workflow Statuses:** Submitted: 1
- Workflow Configuration:** thelagen\_pnl/Terra\_2\_NCBL\_PHB
- Submitted by:** eija.trees@thelagen.cloud, Apr 22, 2024, 7:47 AM
- Total Run Cost:** N/A
- Data Entity:** NCBL\_upload240422, quality\_control\_set
- Submission ID:** 387be0c7-88b6-4c00-8028-9020dc5f08b8
- Call Caching:** Enabled
- Comment:** NCBI upload of 3 sequences to the validation ...
- Delete Intermediate Outputs:** Disabled
- Use Reference Disks:** Disabled
- Retry with More Memory:** Disabled

Below these sections is a table with columns: Data Entity, Last Changed, Status, Run Cost, Messages, and Workflow ID. The table contains one entry for the data entity 'NCBL\_upload240422 (quality\_control\_set)' with a status of 'Queued' and a last changed time of 'Apr 22, 2024, 7:47 AM'.

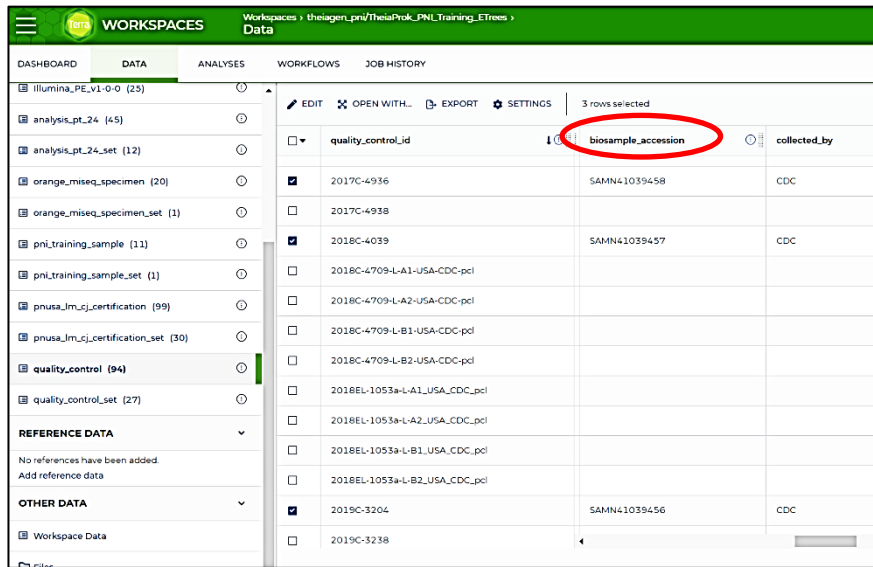
5.7.3.10 Перейдіть на вкладку "Історія завдань", щоб перевірити статус вашого завдання. Успішно виконане завдання позначається зеленою галочкою.

Submission (click for details)	Data entity	No. of Workflows	Status	Submitted	Submission ID	Comment	Actions
Terra_2_NCBL_PHB Submitted by eija.trees@thelagen.cloud	NCBL_upload240422 (qualit...	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 eija.trees@thelagen.cloud	NCBL_upload2024a0117 L...	1	✔ Done	Apr 17, 2024, 7:51 AM	9a8ca79-286a-434a-a684-a38c8f16532	NCBI submission of 3 new sample...	ⓘ
TheiaProk_illumina_PE Submitted by eija.trees@thelagen.cloud	TheiaProk_illumina_PE_202...	6	✔ Done	Apr 15, 2024, 12:52 PM	4b02064-f0b-4942-863c-30830445348	6 strains from the AMD incubator ...	ⓘ
TheiaProk_illumina_PE Submitted by eija.trees@thelagen.cloud	TheiaProk_illumina_PE_202...	3	✔ Done	Apr 15, 2024, 1:12 PM	5d845bc-538a-4d0a-8e3b-66734a07627	Clony pick 2 for 2023 IT strains	ⓘ
TheiaProk_illumina_PE Submitted by eija.trees@thelagen.cloud	TheiaProk_illumina_PE_202...	3	✔ Done	Apr 15, 2024, 12:52 PM	88776c1-9e3f-43a3-a377-a468a47f62a	3 IT isolates from 2023	ⓘ

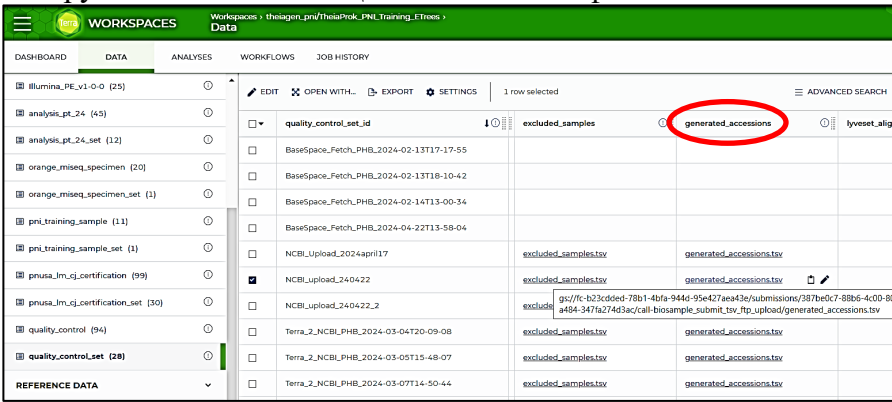
5.7.3.11 Відстеження номерів доступу до NCBI:

5.7.3.11.1 Номери приналежності біозразків (номери SAMN)

5.7.3.11.1.1 У вкладці "Дані" перейдіть до відповідної таблиці даних, і ви побачите поле "biosample\_accession", заповнене номером SAMN, який присвоюється кожному унікальному біозразку.

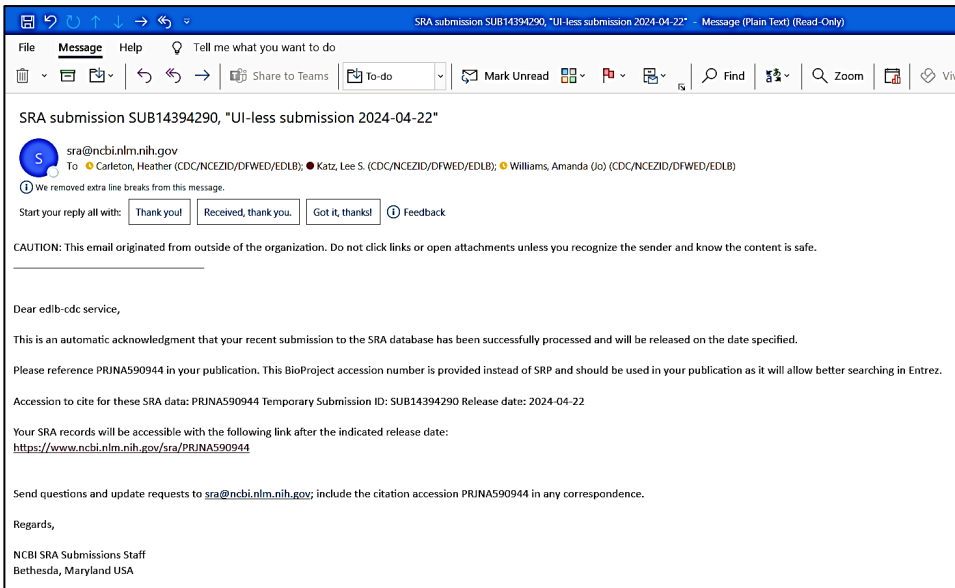


5.7.3.11.1.2 У вкладці "Дані" перейдіть до відповідного набору таблиць даних, знайдіть набір даних, який ви створили для завантаження в NCBI, а потім натисніть на посилання "згенеровані\_доступи" для tsv-файлу, в якому перераховані номери біозразків для завантаженого набору послідовностей. Щоб завантажити файл:

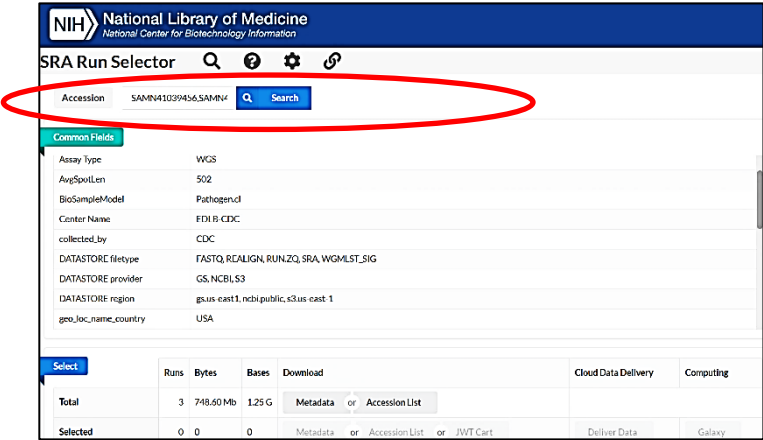


5.7.3.11.1.2.1 Натисніть "Завантажити >\$0.01\*", а потім натисніть "Готово".  
 5.7.3.11.1.2.2 Завантажений файл з'явиться у верхньому правому куті.



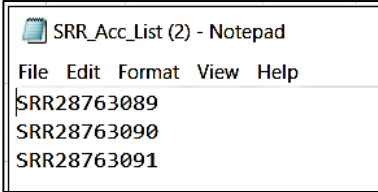
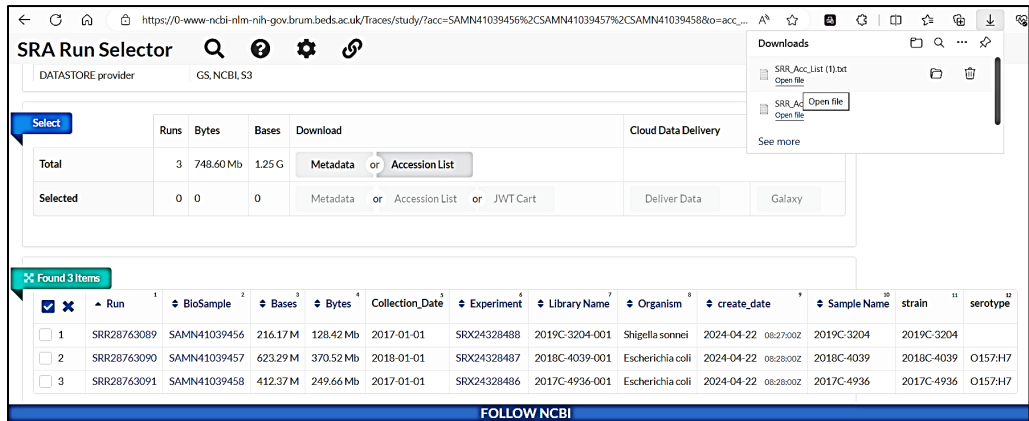


5.7.3.11.2.1 Скопіюйте та вставте номери SAMN з файлу tsv з кроку 5.7.3.11.1.2. або файлу txt з кроку 5.7.3.11.1.3. до інструменту "NCBI Run Selector" за [адресою: https://0-www-ncbi-nlm-nih-gov.brum.beds.ac.uk/Traces/study/](https://0-www-ncbi-nlm-nih-gov.brum.beds.ac.uk/Traces/study/). Розділіть числа комами. Натисніть на кнопку "Пошук".



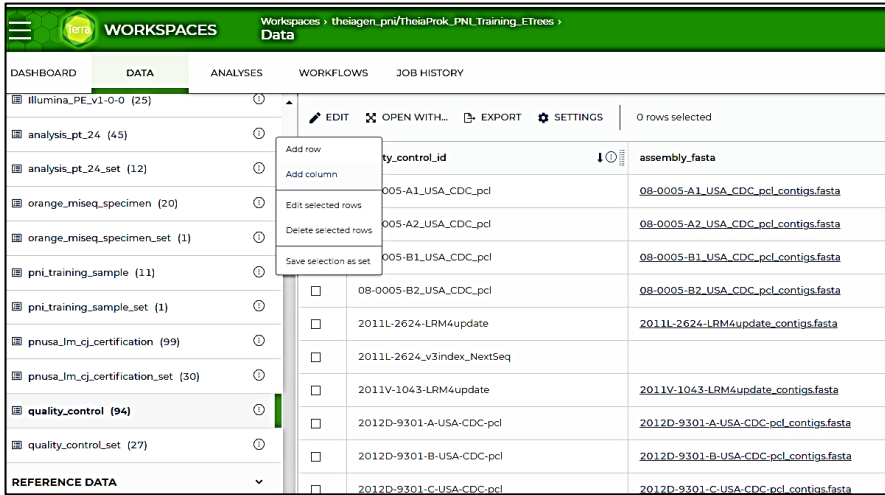
5.7.3.11.2.2 Пошук поверне таблицю, яка містить всю інформацію NCBI про ваші послідовності, включно з номерами доступу до SRR.

5.7.3.11.2.3 Натисніть на "Accession List", щоб завантажити таблицю з номерами присланий. Завантажений файл з'явиться у верхньому правому куті.



5.7.3.11.2.4 Ви можете відстежувати номери приєднання до SRA або в окремій таблиці Excel, або в системі LIMS, або ви можете створити поле "sra\_accession" у вашій таблиці даних Terra і скопіювати та вставити туди номери приєднання для кожного зразка запису. Створити стовпець sra\_accession і використовувати його для відстеження:

5.7.3.11.2.4.1 У випадаючому меню "Редагувати" виберіть "Додати стовпчик".



5.7.3.11.2.4.2 У спливаючому вікні "Додати новий стовпчик" назвіть новий стовпчик "sra\_accession" і натисніть "Зберегти".

**Add a new column**

**Column name**

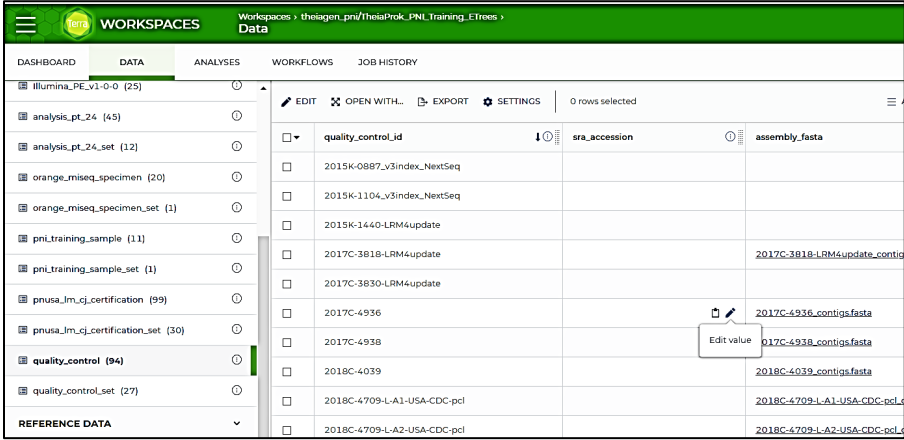
**Default value** (optional, will be entered for all rows)

Type:

String  
  Reference  
  Number  
  Boolean

Value is a list

5.7.3.11.2.4.3 Новий стовпець має з'явитися в таблиці даних. Щоб ввести номер приєднання SRR для певної послідовності, натисніть "Редагувати значення" у стовпчику sra\_accession для цього зразка.



5.7.3.11.2.4.4 У спливаючому вікні "Змінити значення" вставте номер приєднання до SRA у поле та натисніть "Зберегти зміни".

**Edit value**

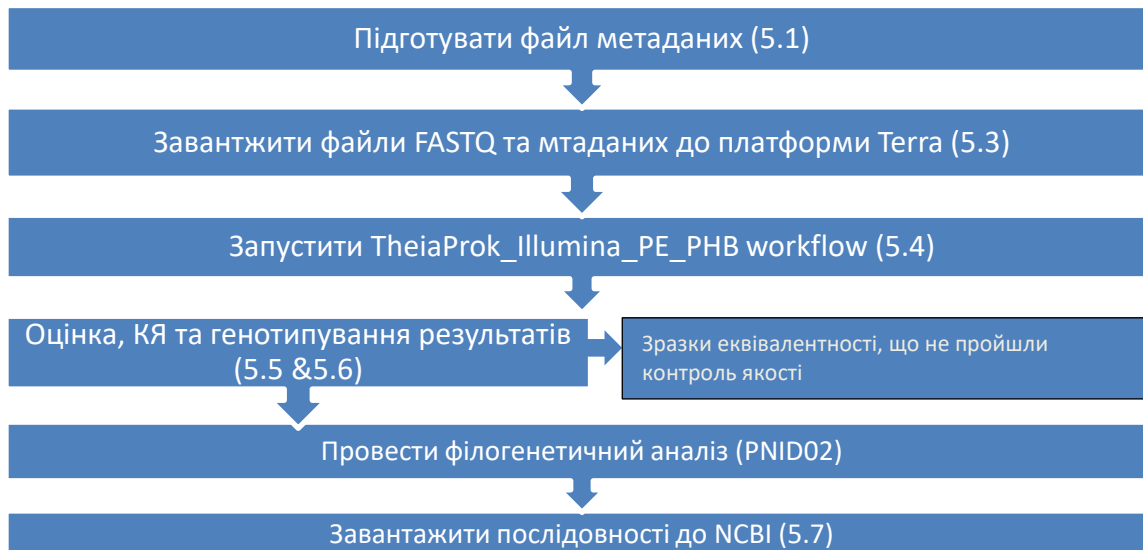
Type:

String  
  Reference  
  Number  
  Boolean

Value is a list

<b>МІЖНАРОДНА СТАНДАРТНА ОПЕРАЦІЙНА ПРОЦЕДУРА PULSENET ДЛЯ АНАЛІЗУ ДАНИХ КОРОТКОГО ЗЧИТУВАННЯ WGS ILLUMINA З ВИКОРИСТАННЯМ ПЛАТФОРМИ TERRA.BI</b>			
Док. No. PNID01	Вер. No. 01	Дата набуття чинності:	Сторінка 31 з 67

## 6. ТЕХНОЛОГІЧНА КАРТА:



## 7. ПОВ'ЯЗАНІ ДОКУМЕНТИ:

7.1 **PNID02:** Міжнародна стандартна операційна процедура PulseNet для філогенетичного аналізу даних WGS з використанням платформи Terra.Bio.

## 8. ВИКОРИСТАНА ЛІТЕРАТУРА:

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. КОНТАКТИ:

9.1 Лабораторія NGS PulseNet Центру контролю та профілактики захворювань США: [pulsenetngslab@cdc.gov](mailto:pulsenetngslab@cdc.gov)

9.2 Міжнародний координатор із забезпечення якості PulseNet Ейя Тріес: [ehyytia-trees@cdc.gov](mailto:ehyytia-trees@cdc.gov)

9.3 Theiagen:

9.3.1 Загальна електронна пошта для підтримки: [support@theiagen.com](mailto:support@theiagen.com)

9.3.2 Мішель Скрібнер: [michelle.scribner@theiagen.com](mailto:michelle.scribner@theiagen.com)

9.3.3 Френк Амбросіо: [frank.ambrosio@theiagen.com](mailto:frank.ambrosio@theiagen.com)

10. ЗМІНИ: Відсутні.

<b>МІЖНАРОДНА СТАНДАРТНА ОПЕРАЦІЙНА ПРОЦЕДУРА PULSENET ДЛЯ АНАЛІЗУ ДАНИХ КОРОТКОГО ЗЧИТУВАННЯ WGS ILLUMINA З ВИКОРИСТАННЯМ ПЛАТФОРМИ TERRA.VI</b>			
Док. No. PNID01	Вер. No. 01	Дата набуття чинності:	Сторінка 32 з 67

### 11. ПІДПИСИ ПРО ЗГОДУ:

Затверджено \_\_\_\_\_ Дата: \_\_\_\_\_  
Персонал PulseNet ЗЯ/КЯ

Затверджено \_\_\_\_\_ Дата: \_\_\_\_\_  
Технічний керівник PulseNet WGS

Затверджено \_\_\_\_\_ Дата: \_\_\_\_\_  
Міжнародний координатор PulseNet

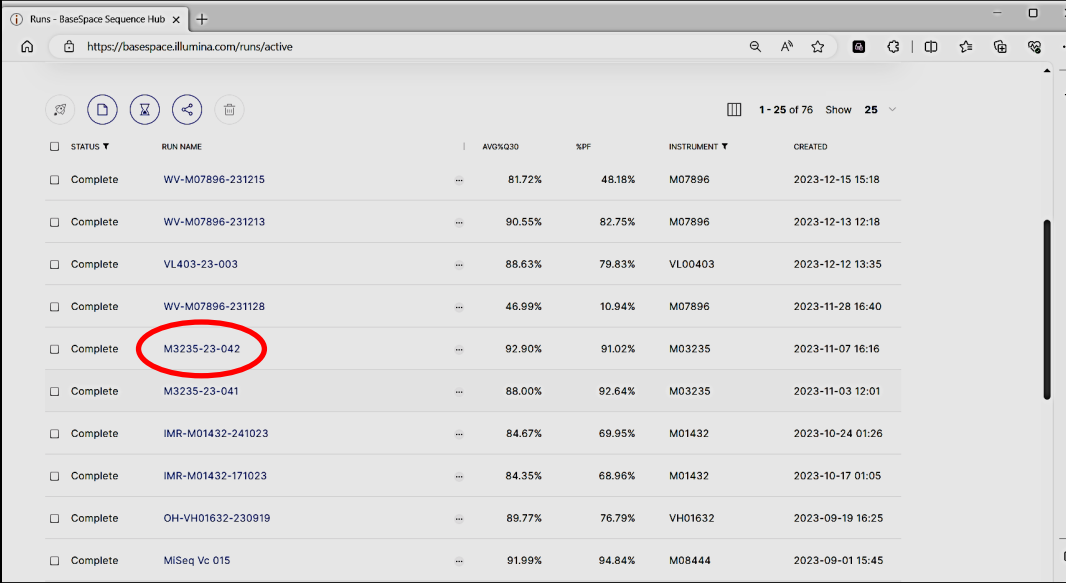
Затверджено \_\_\_\_\_ Дата: \_\_\_\_\_  
Керівник групи реагування та управління спалахами PulseNet

Затверджено \_\_\_\_\_ Дата: \_\_\_\_\_  
Завідувач сектору лабораторії кишкових захворювань

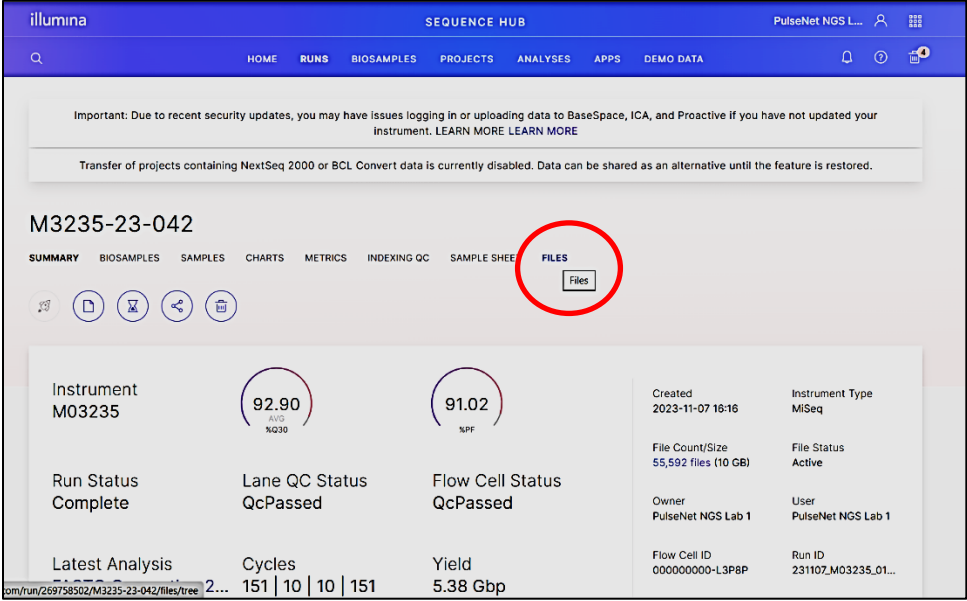
**Додаток PNID01-1: Імпорт даних до Terra безпосередньо з Illumina BaseSpace**

**ПРИМІТКА:** Щоб налаштувати робочу область Terra для підключення до облікового запису Illumina BaseSpace, дотримуйтеся інструкцій на сайті ресурсів Theiagen за адресою [https://theiagen.notion.site/BaseSpace\\_Fetch-34978656aa2d46ba82f2059434bd9369](https://theiagen.notion.site/BaseSpace_Fetch-34978656aa2d46ba82f2059434bd9369). Для отримання додаткової допомоги зверніться до компанії Theiagen (контактну інформацію див. у розділі *Контакти (9)*).

1. Увійдіть у свій обліковий запис BaseSpace і знайдіть цикл, який потрібно імпортувати до Terra.



2. Завантажте Зразок для пробігу:
  - a. Перейдіть на вкладку "Файли" і прокрутіть сторінку донизу.



- b. Натисніть на посилання SampleSheet.csv.

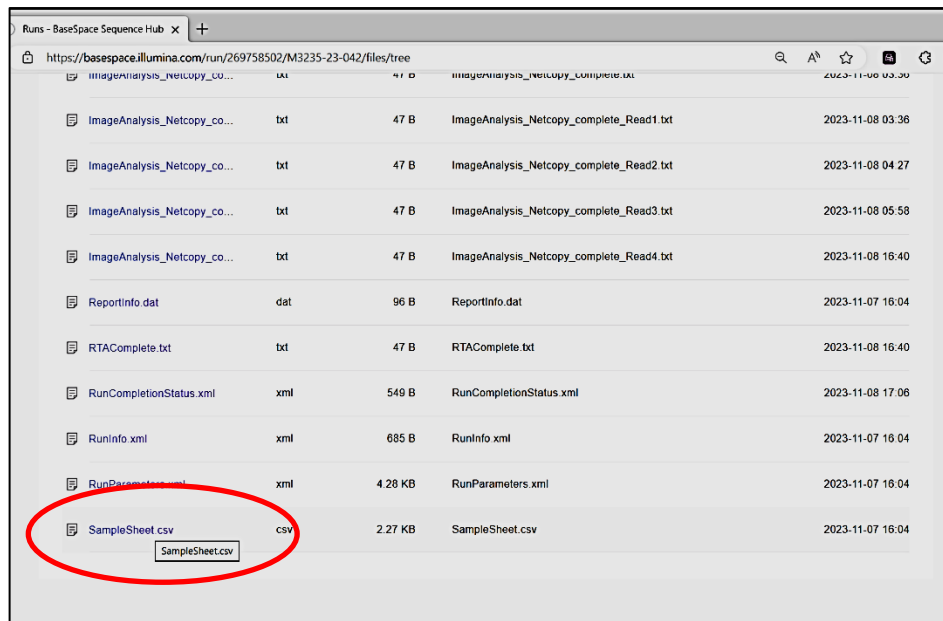
**МІЖНАРОДНА СТАНДАРТНА ОПЕРАЦІЙНА ПРОЦЕДУРА PULSENET ДЛЯ АНАЛІЗУ ДАНИХ  
КОРОТКОГО ЗЧИТУВАННЯ WGS ILLUMINA З ВИКОРИСТАННЯМ ПЛАТФОРМИ TERRA.VI**

Док. No. PNID01

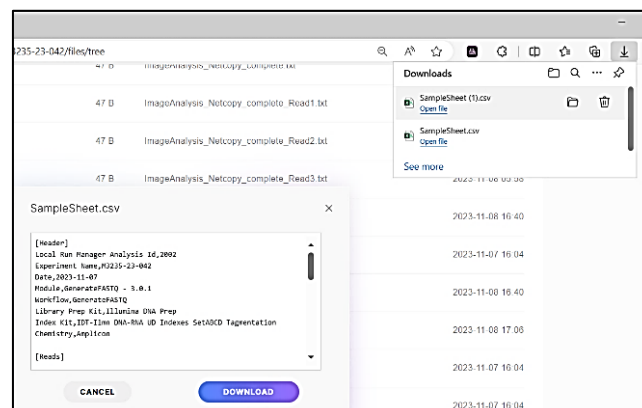
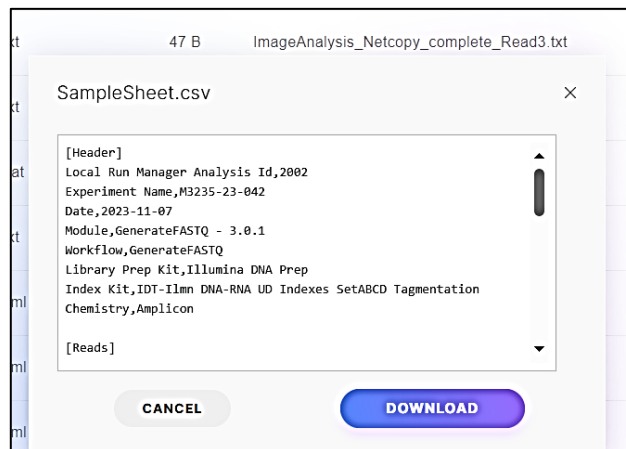
Вер. No. 01

Дата набуття чинності:

Сторінка 34 з 67



с. У спливаючому вікні "SampleSheet.csv" натисніть "Завантажити".



- d. Завантажений csv-файл з'явиться у верхньому правому куті в розділі "Завантаження".
3. Відкрийте аркуш SampleSheet. Стовпці, необхідні у файлі метаданих tsv, залежатимуть від стовпців і вмісту стовпців, наявних на аркуші SampleSheet.
4. Підготуйте файл метаданих tsv:
  - a. Стовпці потрібні, коли стовпці "Sample\_Name" та "Sample\_ID" в аркуші SampleSheet мають **однаковий** вміст:

Sample_ID	Sample_Name	Description	Index_Plat	Index_Plat	I7_Index	I_Index	I5_Index	Index2	Sample_Project
05480-M3235-23-042	D5480-M3235-23-042		B	A06	UDP0137	CCGGTTC	UDP0137	TATATTCGAG	
ATCC-BAA-460-M3235-23-042	ATCC-BAA-460-M3235-23-042		B	B06	UDP0138	GGCCAAT	UDP0138	CGGTCCGATA	
2011L-2624-M3235-23-042	2011L-2624-M3235-23-042		B	C06	UDP0139	GAATACCT	UDP0139	ACAATAGAGT	
2015K-0092-M3235-23-042	2015K-0092-M3235-23-042		B	D06	UDP0140	TACGTGA	UDP0140	CGGTATTAG	
2015K-1440-M3235-23-042	2015K-1440-M3235-23-042		B	E06	UDP0141	CITATTGG	UDP0141	GATAACAAGT	
2013V-1178-M3235-23-042	2013V-1178-M3235-23-042		B	F06	UDP0142	ACAATACT	UDP0142	AGTTATCACA	
2014C-3598-M3235-23-042	2014C-3598-M3235-23-042		B	G06	UDP0143	GTTGGATC	UDP0143	TTCCAGGTAA	
2014C-3857-M3235-23-042	2014C-3857-M3235-23-042		B	H06	UDP0144	AATCCAAT	UDP0144	CATGTAGAGG	
2015C-3881-M3235-23-042	2015C-3881-M3235-23-042		B	A07	UDP0145	TATGATGC	UDP0145	GATTGTGATA	

- i. ідентифікатор\_імені\_бази\_даних\_об'єкта:
  1. Введіть назву таблиці даних (нової або існуючої) в комірку A1 між "entity:" та "id".
  2. Введіть ідентифікатори зразків у стовпчик A так, як ви хочете, щоб вони з'явилися в таблиці даних Terra.
- ii. basespace\_sample\_name: скопіюйте та вставте вміст з поля SampleSheet "Sample\_Name".
- iii. basespace\_collection\_id: введіть назву циклу так, як вона відображається у BaseSpace.

	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. Стівці, необхідні, коли стівці "Sample\_Name" та "Sample\_ID" на аркуші SampleSheet мають **різний** вміст:

	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-iIrn DNA-RNA UD Indexes SetABCD Tagmentation											
9	Chemistry	Amplicon											
10													
11	[Reads]												
12		151											
13		151											
14													
15	[Settings]												
16	adapter	CTGTCTCTATACACATCT											
17													
18	[Data]												
19	Sample_ID	Sample_Name	Descriptio	Index	Pls	I7_Index	index	I5_Index	Index2	Sample_Project			
20	2023FD-00134	BE230960535-CAOC-M5870-230530E	A	A07	UDP0049	AGTGTG/UDP0049	CTGGTAC/CPD_230530E						
21	BE230960535	BE230960535-CAOC-M5870-230530E	A	B07	UDP0050	GACACCA/UDP0050	TCAACGT/Sal_230530E						
22	2023FD-00135	2023FD-00135-CAOC-M5870-230530E	A	C07	UDP0051	CCTGTCT/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	TTCGGTCE/Hiq230530E						
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	CGTCCAG/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. Введіть назву таблиці даних (нової або існуючої) в комірку A1 між "entity:" та "id".
  2. Введіть ідентифікатори зразків у стовпчик A так, як ви хочете, щоб вони з'явилися в таблиці даних Terra.

- ii. `basespace_sample_name`: скопіюйте та вставте вміст з поля SampleSheet "Sample\_Name".
- iii. `basespace_sample_id`: скопіюйте та вставте вміст з поля SampleSheet "Sample\_ID".
- iv. `basespace_collection_id`: введіть назву циклу так, як вона відображається у 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									

**с. Стовпці, необхідні для NextSeq SampleSheet:**

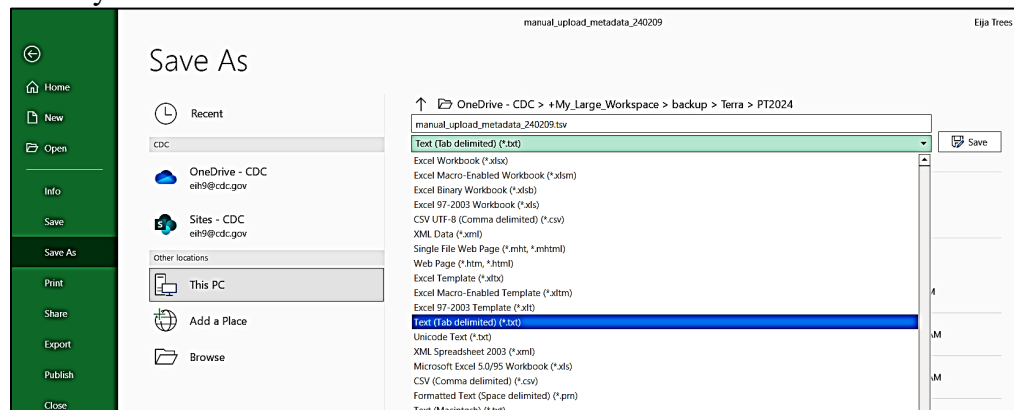
	A	B	C
1	[Header]		
2	FileFormatVersion		2
3	RunName	VL403-24-001	
4	InstrumentPlatform	NextSeq1k2k	
5	IndexOrientation	Forward	
6			
7	[Reads]		
8	Read1Cycles		151
9	Read2Cycles		151
10	Index1Cycles		10
11	Index2Cycles		10
12			
13	[Sequencing_Settings]		
14	LibraryPrepKits	illuminaDNAPrep	
15			
16	[BCLConvert_Settings]		
17	SoftwareVersion	3.10.12	
18	AdapterRead1	CTGTCTCTTATACACATCT	
19	AdapterRead2	CTGTCTCTTATACACATCT	
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	TAGTTCGGTA	CCATGTGTAG
29	2013L-5585	CTATTACTAC	GAGTCTCC
30	2013L-5615	TAGCATAACC	GCTATGGCCA
31	2011L-2624	ACTCTATTGT	ATCGCATATG
32	2015K-9887	CCAAAGGCTT	TCGAAGTACT
33	2015K-1104	TTACTCCACA	GACACCGATG
34	2014K-0833	AGTAGAAGTG	CTAGCGTCGA

- i. `entity_datablename_id`:
  1. Введіть назву таблиці даних (нової або існуючої) в комірку A1 між "entity:" та "id".
  2. Введіть ідентифікатори зразків у стовпчик A так, як ви хочете, щоб вони з'явилися в таблиці даних Terra.
- ii. `basespace_sample_name`: скопіюйте та вставте вміст з поля SampleSheet "Sample\_ID".
- iii. `basespace_sample_id`: скопіюйте та вставте вміст з поля SampleSheet "Sample\_ID".

- iv. `basespace_collection_id`: введіть назву циклу так, як вона відображається у BaseSpace.

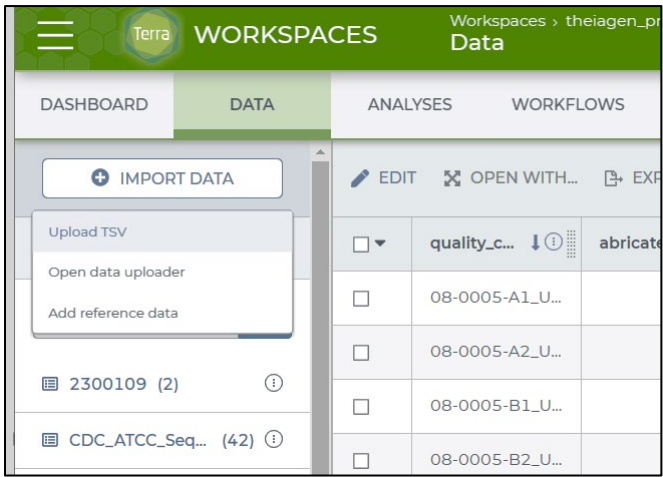
entity_quality_control_id	basespace_sample_name	basespace_sample_id	basespace_collection_id
2013L-5214_v3index_NextSeq	2013L-5214	2013L-5214	VL403-24-001
2013L-5351_v3index_NextSeq	2013L-5351	2013L-5351	VL403-24-001
2013L-5356_v3index_NextSeq	2013L-5356	2013L-5356	VL403-24-001
2013L-5357_v3index_NextSeq	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. Збережіть файл у **форматі tsv**: виберіть "Зберегти як" і "Текст (з розділенням табуляцією) (\*.txt)". Переконайтеся, що ім'я вашого файлу закінчується на **".tsv"**.

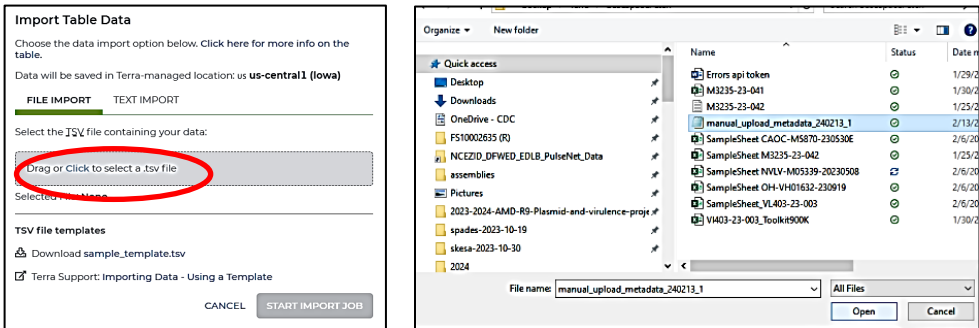


5. Імпортуйте файл метаданих tsv:

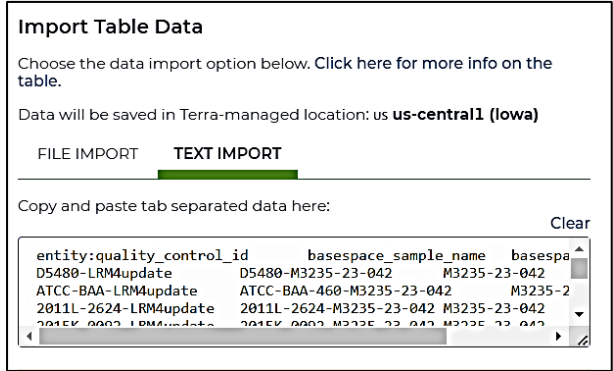
- a. На вкладці "Дані" натисніть "Імпортувати дані" і виберіть "Завантажити tsv" у випадаючому меню.



- b. У спливаючому вікні "Імпортувати дані таблиці" ви можете імпортувати метадані двома способами:
- i. На вкладці "Імпорт файлів" натисніть посередині, щоб вибрати файл tsv, перейдіть до місця, де зберігається файл метаданих tsv, виберіть файл і натисніть "Відкрити".

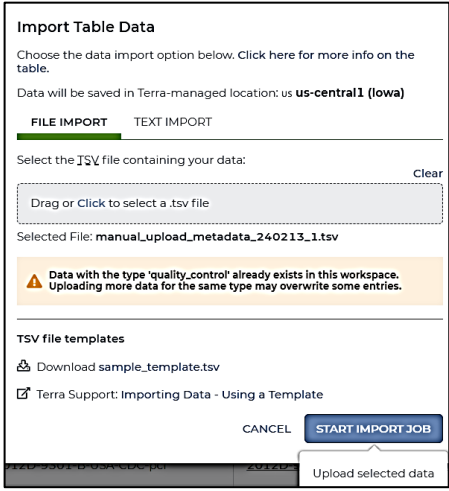


- ii. Крім того, ви можете перейти на вкладку "Імпорт тексту" і скопіювати та вставити вміст tsv-файлу в поле посередині.

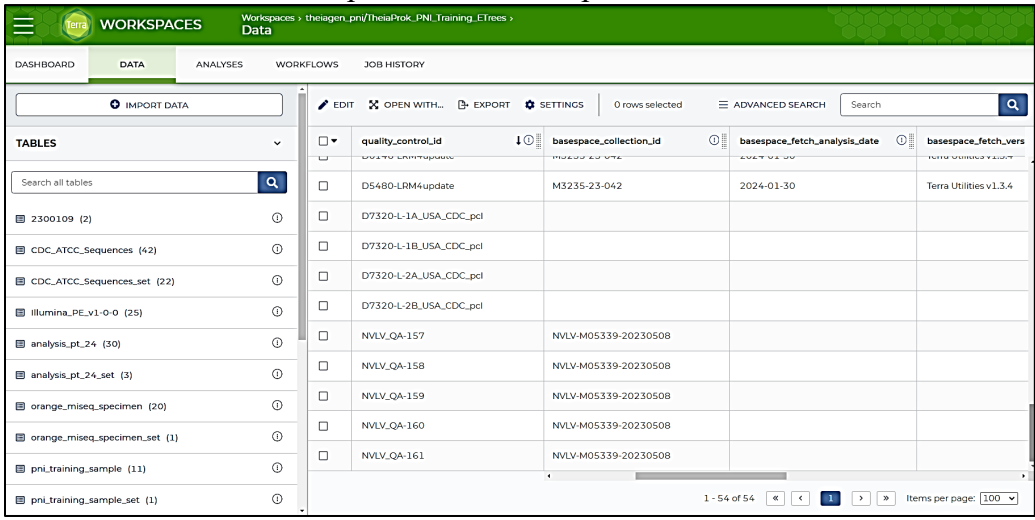


- c. У спливаючому вікні "Імпортувати дані з таблиці" з'явиться попередження про те, що дані вже існують у відповідній таблиці даних (якщо імпорт

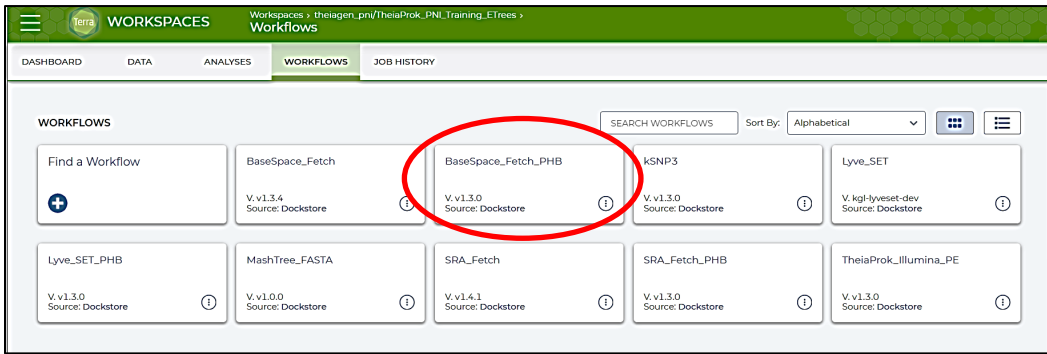
виконується в існуючу таблицю даних), і завантаження нових даних до неї може призвести до перезапису існуючих даних. Натисніть "Почати імпорт".



d. Після завершення імпорту ви побачите нові записи, створені в таблиці даних для послідовностей, імпортованих з BaseSpace.



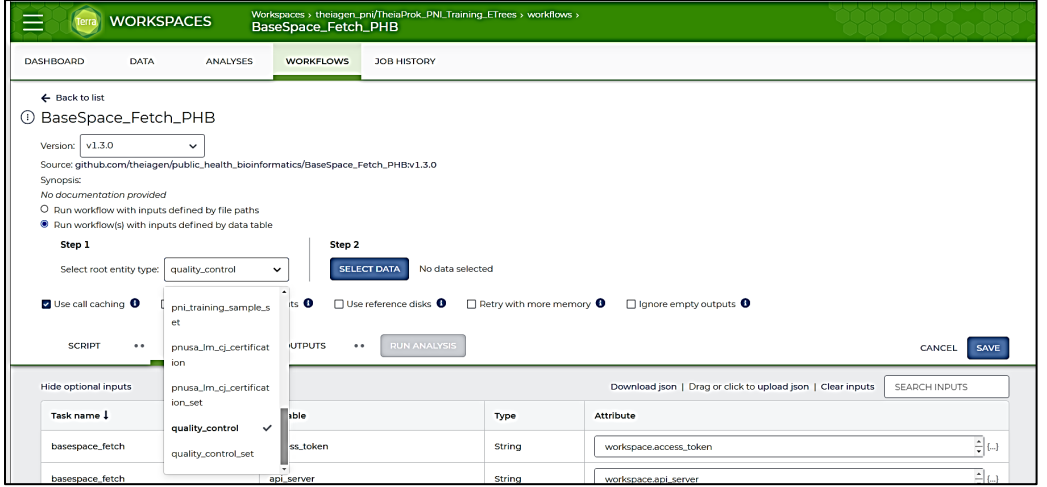
6. Запустіть робочий процес "BaseSpace\_Fetch\_PNB":  
 а. На вкладці "Робочі процеси" натисніть "BaseSpace\_Fetch\_PNB". Відкриється вікно "BaseSpace\_Fetch\_PNB".



b. У випадяючому меню "Версія" виберіть останню версію BaseSpace\_Fetch\_PHB.



c. У розділі "Крок 1" натисніть на випадяюче меню "Select root entity type" і виберіть таблицю даних, до якої ви завантажили (кроки 4-5) tsv-файл, що містить метадані для циклу, який потрібно імпортувати з BaseSpace.



d. У розділі "Крок 2" натисніть "Вибрати дані" (скріншот вище). Ви потрапите на екран вибору вибірки.

e. Встановіть прапорці навпроти зразків, які потрібно імпортувати з BaseSpace, і натисніть "ОК". Ви повернетесь до вікна "BaseSpace\_Fetch\_PHB".

**МІЖНАРОДНА СТАНДАРТНА ОПЕРАЦІЙНА ПРОЦЕДУРА PULSENET ДЛЯ АНАЛІЗУ ДАНИХ  
КОРОТКОГО ЗЧИТУВАННЯ WGS ILLUMINA З ВИКОРИСТАННЯМ ПЛАТФОРМИ TERRA.VI**

Док. No. PNID01

Вер. No. 01

Дата набуття чинності:

Сторінка 42 з 67

Select Data

Choose specific quality\_controls to process  
Choose existing sets of quality\_controls

Select quality\_controls to process SETTINGS 5 rows selected

quality_control_id	assembly_fasta	basespace_collection_id	basespace_fetch_analysis_date	basespace_fetch_version
<input type="checkbox"/>	D7320-L-2B_USA_CDC_pcl	D7320-L-2B_USA_CDC_pcl_contigs.fasta		
<input checked="" type="checkbox"/>	NVLV_QA-157		NVLV-M05339-20230508	
<input checked="" type="checkbox"/>	NVLV_QA-158		NVLV-M05339-20230508	
<input checked="" type="checkbox"/>	NVLV_QA-159		NVLV-M05339-20230508	
<input checked="" type="checkbox"/>	NVLV_QA-160		NVLV-M05339-20230508	
<input checked="" type="checkbox"/>	NVLV_QA-161		NVLV-M05339-20230508	

Selected quality\_controls will be saved as a new quality\_control\_set named:  
BaseSpace\_Fetch\_PHB\_2024-02-13T17-17-55

CANCEL OK

f. Зніміть прапорець з пункту "Use call caching".

No documentation provided

Run workflow with inputs defined by file paths  
Run workflow(s) with inputs defined by data table

Step 1  
Select data table: quality\_control

Step 2  
5 selected quality\_controls (will create a new quality\_control\_set named "BaseSpace\_Fetch\_PHB\_2024-05-23T17-14-14")

Use call caching  Delete intermediate outputs  Use reference disks  Retry with more memory  Ignore empty outputs

SCRIPT \*\* INPUTS \*\* OUTPUTS \*\* Run Analysis

Hide optional inputs Download json | Drag or click

Task name	Variable	Type	Input value
-----------	----------	------	-------------

g. На вкладці "Inputs" визначте наступні змінні в полі "Атрибут":

**ПРИМІТКА:** Коли ви заповнюєте стовпчик Атрибут, клацання всередині клітинки відкриє спадне меню атрибутів, які ви можете вибрати, щоб уникнути помилок.

- i. access\_token: "workspace.access\_token".
- 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.dataslename\_id", наприклад, "this.quality\_control\_id".
- vi. basespace\_sample\_id: "this.basespace\_sample\_id".

**ПРИМІТКА:** потрібно заповнювати лише в тому випадку, якщо вміст полів "Sample\_Name" та "Sample\_ID" відрізняються у файлі SampleSheet циклу або ви імпортуєте дані з циклу NextSeq.

**МІЖНАРОДНА СТАНДАРТНА ОПЕРАЦІЙНА ПРОЦЕДУРА PULSENET ДЛЯ АНАЛІЗУ ДАНИХ  
КОРОТКОГО ЗЧИТУВАННЯ WGS ILLUMINA З ВИКОРИСТАННЯМ ПЛАТФОРМИ TERRA.VI**

Док. No. PNID01

Вер. No. 01

Дата набуття чинності:

Сторінка 43 з 67

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. На вкладці "Outputs" натисніть "Використовувати за замовчуванням", потім "Зберегти", а потім "Запустити аналіз".

**ПРИМІТКА:** Кнопка "Зберегти" відображається лише в тому випадку, якщо параметри (крім ідентифікаторів зразків) були змінені з моменту попереднього надсилання завдання. Кнопка "Виконати аналіз" стає підсвіченою після збереження параметрів.

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. У спливаючому вікні "Підтвердити запуск" опишіть вашу заявку (необов'язково) і натисніть "Запустити".

**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. Відкриється вкладка "Історія робіт", де статус поданих заявок спочатку має бути "В черзі".

The screenshot shows the Terra Workspaces interface for a specific job history entry. The workflow is in a 'Queued' state. The configuration includes a workflow named 'theiaigen\_pnl/BaseSpace\_Fetch\_PHB' and a data entity 'BaseSpace\_Fetch\_PHB\_2024-02-13T17-17-55\_quality\_control\_set'. The submission was made by 'eja:rees@theiaigen.cloud' on Feb 13, 2024, at 12:25 PM. The total run cost is N/A. The workflow configuration includes options for 'Delete Intermediate Outputs' (Disabled), 'Use Reference Disks' (Disabled), and 'Retry with More Memory' (Disabled). The 'Call Caching' is Enabled. A comment is present: 'NVLV MISEq Basespace run. Sample\_na...'. The 'WORKFLOWS' tab is active, showing a table of workflow instances.

Data Entity ↓	Last Changed	Status	Run Cost	Message	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		

к. Після завершення роботи буде показано статус "Успішно" або "Виконано". На вкладці "Дані" ви побачите імена файлів FASTQ у стовпчиках "Read1" та "Read2", а також інформацію у стовпчиках "basespace\_fetch\_analysis\_date" і "basespace\_fetch\_version".

The screenshot shows the Terra Workspaces interface for the same job history entry, but now the workflow is in a 'Succeeded' state. The configuration and submission details are the same as in the previous screenshot. The 'Workflow Statuses' section now shows a green checkmark and 'Succeeded: 5'. The 'WORKFLOWS' tab is active, showing a table of workflow instances, all of which are now 'Succeeded'.

Data Entity ↓	Last Changed	Status	Run Cost	Message	Workflow ID	Links
NVLV_QA-157 (quality_control)	Feb 13, 2024, 12:30 PM	✔ Succeeded	N/A		6154609d-c6dc-4d1d-e1c8-5286b71303...	🔗 📄 🗑️
NVLV_QA-158 (quality_control)	Feb 13, 2024, 12:29 PM	✔ Succeeded	N/A		d3535af4-758a-490d-9be9-166964578000	🔗 📄 🗑️
NVLV_QA-159 (quality_control)	Feb 13, 2024, 12:29 PM	✔ Succeeded	N/A		1e1351e2-5f82-4978-9777-0ecba6e59472	🔗 📄 🗑️
NVLV_QA-160 (quality_control)	Feb 13, 2024, 12:29 PM	✔ Succeeded	N/A		f561f71f-6ada-455c-8820-8f6fb54b4094	🔗 📄 🗑️
NVLV_QA-161 (quality_control)	Feb 13, 2024, 12:29 PM	✔ Succeeded	N/A		3be0f9ce-a2e3-4302-af5a-01187ac5eb07	🔗 📄 🗑️

**МІЖНАРОДНА СТАНДАРТНА ОПЕРАЦІЙНА ПРОЦЕДУРА PULSENET ДЛЯ АНАЛІЗУ ДАНИХ КОРОТКОГО ЗЧИТУВАННЯ WGS ILLUMINA З ВИКОРИСТАННЯМ ПЛАТФОРМИ TERRA.VI**

Док. No. PNID01

Вер. No. 01

Дата набуття чинності:

Сторінка 45 з 67

WORKSPACES Workspaces > theisgen\_pn/TheisProk\_PNL\_Training\_ETrees > Data

DASHBOARD DATA ANALYSES WORKFLOWS JOB HISTORY

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

	quality_controlId	sample_name	read1	read2
<input type="checkbox"/>	D5480-LRM4update	13235-23-042	D5480-LRM4update_R1.fasta.gz	D5480-LRM4update_R2.fasta.gz
<input type="checkbox"/>	D7320-L-1A_USA_CDC_pcl		D7320-L-1A-M3235-21-007-512_L001_R1_001.fast-	D7320-L-1A-M3235-21-007-512_L001_L-
<input type="checkbox"/>	D7320-L-1B_USA_CDC_pcl		D7320-L-1B-M3235-21-007-513_L001_R1_001.fast-	D7320-L-1B-M3235-21-007-513_L001_L-
<input type="checkbox"/>	D7320-L-2A_USA_CDC_pcl		D7320-L-2A-M3235-21-007-514_L001_R1_001.fast-	D7320-L-2A-M3235-21-007-514_L001_L-
<input type="checkbox"/>	D7320-L-2B_USA_CDC_pcl		D7320-L-2B-M3235-21-007-515_L001_R1_001.fast-	D7320-L-2B-M3235-21-007-515_L001_L-
<input type="checkbox"/>	NVLV_QA-157		NVLV_QA-157_R1.fasta.gz	NVLV_QA-157_R2.fasta.gz
<input type="checkbox"/>	NVLV_QA-158		NVLV_QA-158_R1.fasta.gz	NVLV_QA-158_R2.fasta.gz
<input type="checkbox"/>	NVLV_QA-159		NVLV_QA-159_R1.fasta.gz	NVLV_QA-159_R2.fasta.gz
<input type="checkbox"/>	NVLV_QA-160		NVLV_QA-160_R1.fasta.gz	NVLV_QA-160_R2.fasta.gz
<input type="checkbox"/>	NVLV_QA-161		NVLV_QA-161_R1.fasta.gz	NVLV_QA-161_R2.fasta.gz

1 - 54 of 54 Items per page: 100

WORKSPACES Workspaces > theisgen\_pn/TheisProk\_PNL\_Training\_ETrees > Data

DASHBOARD DATA ANALYSES WORKFLOWS JOB HISTORY

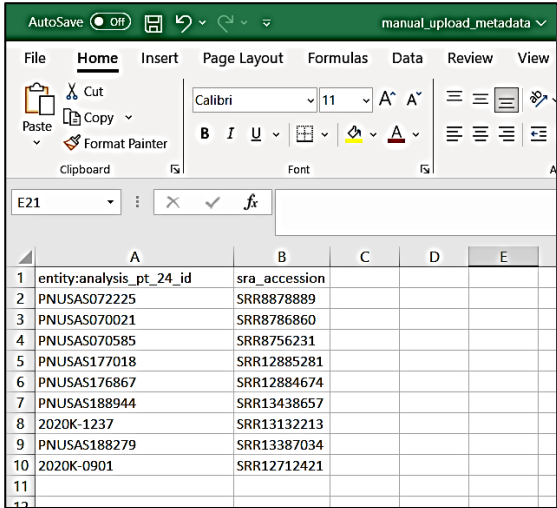
EDIT OPEN WITH... EXPORT SETTINGS 0 rows selected ADVANCED SEARCH Search

	quality_controlId	basespace_fetch_analysis_date	basespace_fetch_version	basespace_sample_name
<input type="checkbox"/>	D5480-LRM4update	2024-01-30	Terra Utilities v1.3.4	D5480-M3235-23-042
<input type="checkbox"/>	D7320-L-1A_USA_CDC_pcl			
<input type="checkbox"/>	D7320-L-1B_USA_CDC_pcl			
<input type="checkbox"/>	D7320-L-2A_USA_CDC_pcl			
<input type="checkbox"/>	D7320-L-2B_USA_CDC_pcl			
<input type="checkbox"/>	NVLV_QA-157	2024-02-13	PHB v1.3.0	QA-157
<input type="checkbox"/>	NVLV_QA-158	2024-02-13	PHB v1.3.0	QA-158
<input type="checkbox"/>	NVLV_QA-159	2024-02-13	PHB v1.3.0	QA-159
<input type="checkbox"/>	NVLV_QA-160	2024-02-13	PHB v1.3.0	QA-160
<input type="checkbox"/>	NVLV_QA-161	2024-02-13	PHB v1.3.0	QA-161

1 - 54 of 54 Items per page: 100

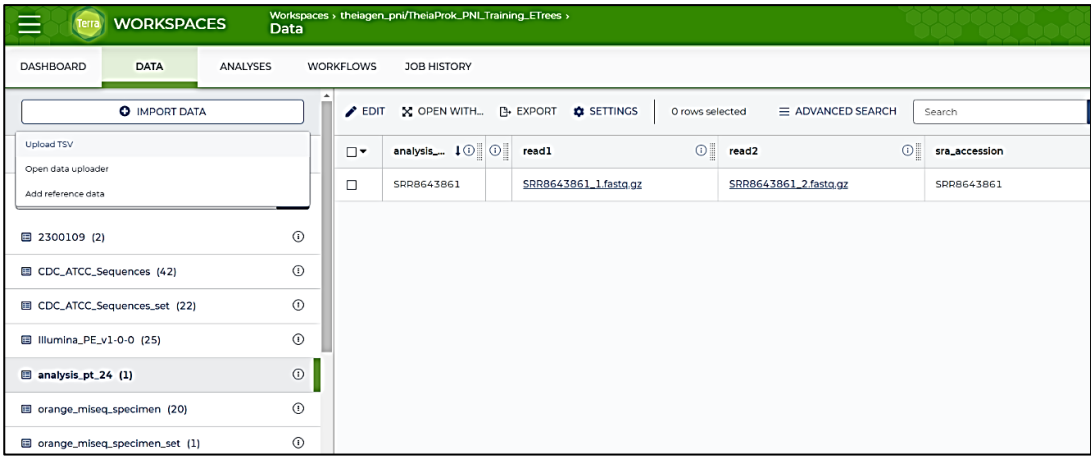
**Додаток PNID01-2: Завантаження даних з NCBI SRA**

1. Підготуйте файл метаданих tsv:
  - a. Введіть назву таблиці даних (нової або існуючої) в комірку A1 між "entity:" та "id".
  - b. Введіть ідентифікатори зразків у стовпчик A так, як ви хочете, щоб вони з'явилися в таблиці даних Terra.
  - c. sra\_accession: введіть номери доступу до SRA для послідовностей, які потрібно завантажити з SRA.

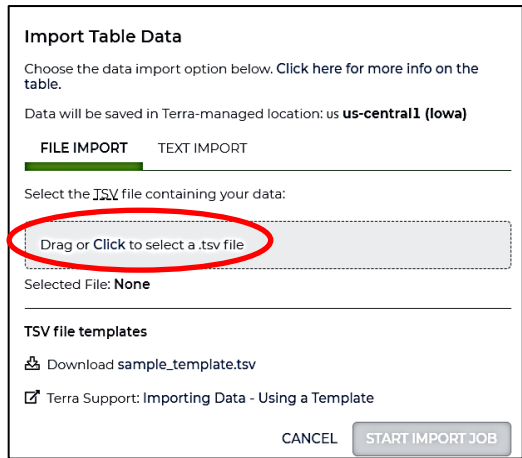


- d. Збережіть файл у форматі tsv: виберіть "Зберегти як" і "Текст (з розділенням табуляцією) (\*.txt)". Переконайтеся, що ім'я вашого файлу закінчується на ".tsv".

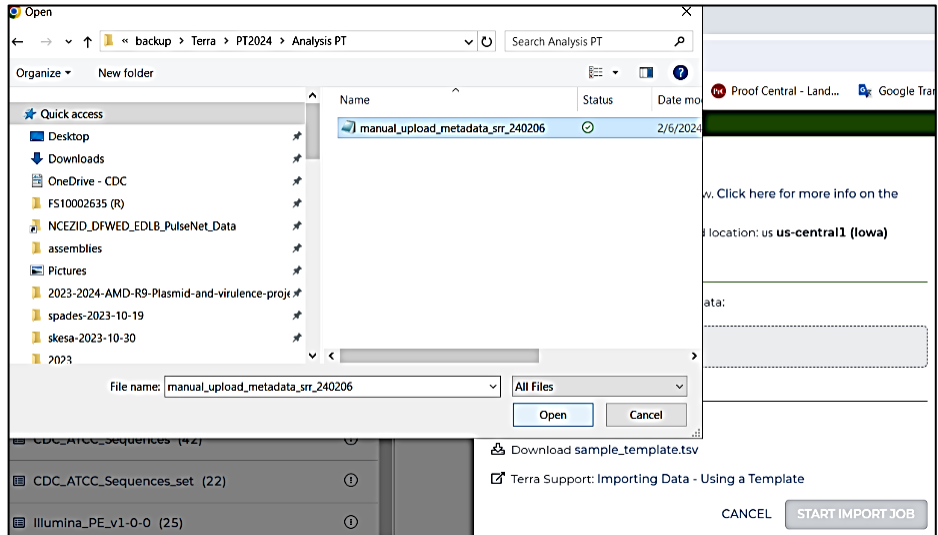
2. На вкладці "Дані" натисніть "Імпортувати дані" і виберіть "Завантажити tsv" у випадаючому меню.



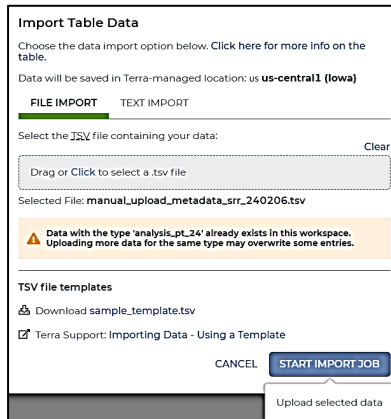
3. У спливаючому вікні "Імпорт даних таблиці" на вкладці "Імпорт файлів" клацніть посередині, щоб вибрати файл tsv.



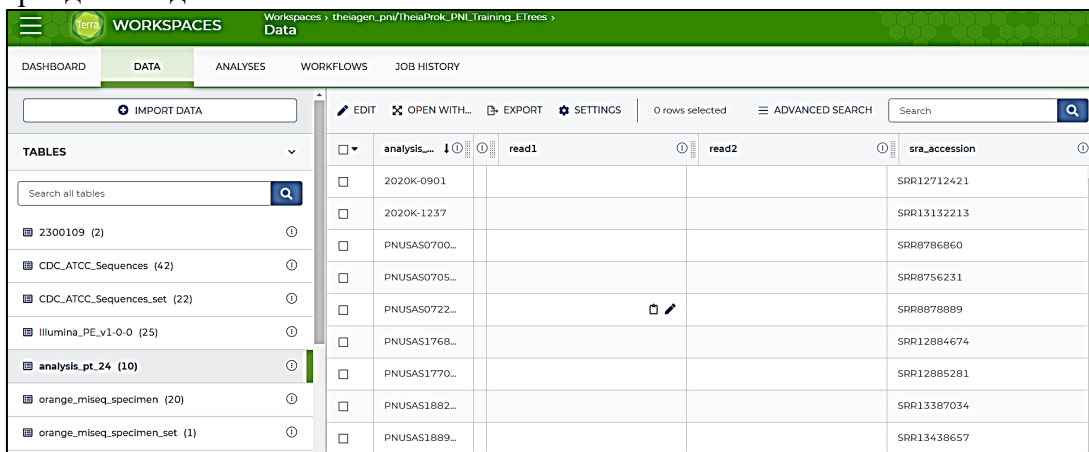
4. Перейдіть до місця, де зберігається файл метаданих tsv, виберіть його і натисніть "Відкрити".



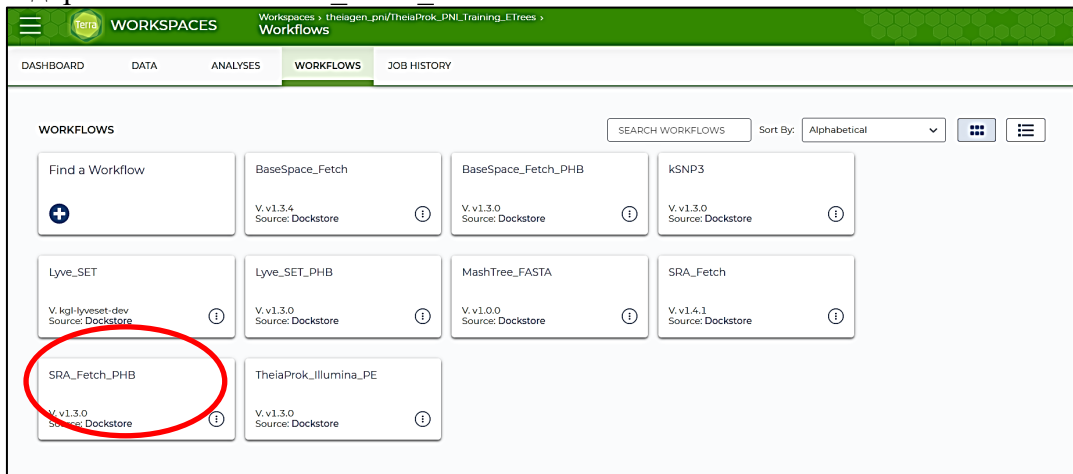
5. У спливаючому вікні "Імпортувати дані з таблиці" з'явиться попередження про те, що у відповідній таблиці даних вже існують дані (якщо імпорт виконується в існуючу таблицю даних), і завантаження нових даних до неї може призвести до перезапису існуючих даних. Натисніть "Почати імпорт".



6. Після завершення імпорту ви побачите нові записи, створені в таблиці даних для послідовностей, які будуть завантажені з NCBI, разом з їхніми номерами приєднання до SRA.



7. На вкладці "Робочі процеси" натисніть на робочий процес "SRA\_Fetch\_PHB". Відкриється вікно "SRA\_Fetch\_PHB".



8. У випадяючому меню "Версія" виберіть останню версію SRA\_Fetch\_PHB.

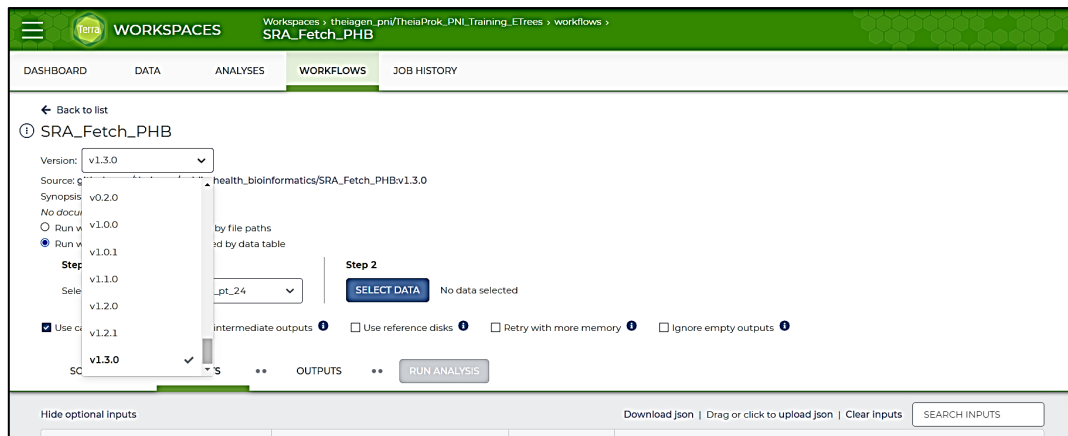
**МІЖНАРОДНА СТАНДАРТНА ОПЕРАЦІЙНА ПРОЦЕДУРА PULSENET ДЛЯ АНАЛІЗУ ДАНИХ  
КОРОТКОГО ЗЧИТУВАННЯ WGS ILLUMINA З ВИКОРИСТАННЯМ ПЛАТФОРМИ TERRA.VI**

Док. No. PNID01

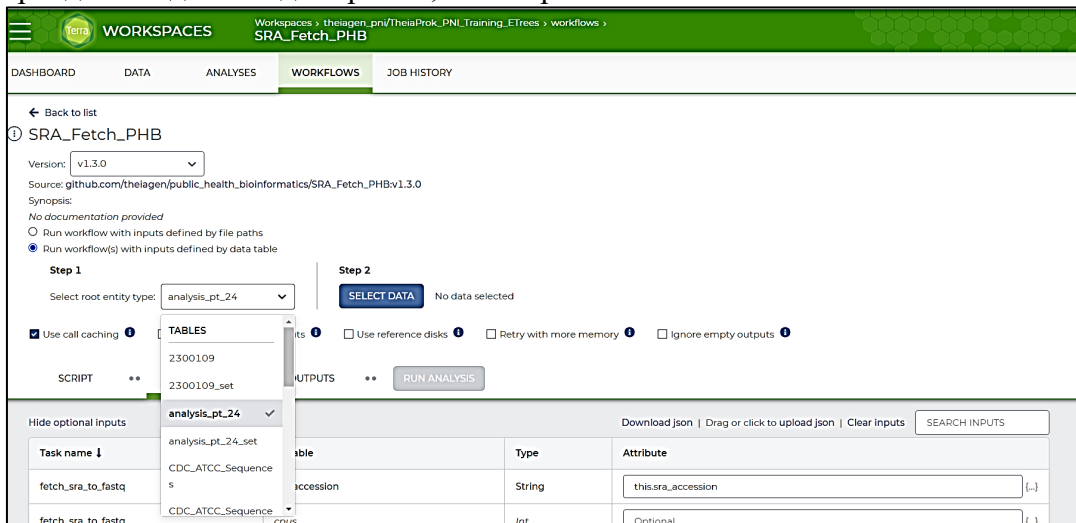
Вер. No. 01

Дата набуття чинності:

Сторінка 49 з 67

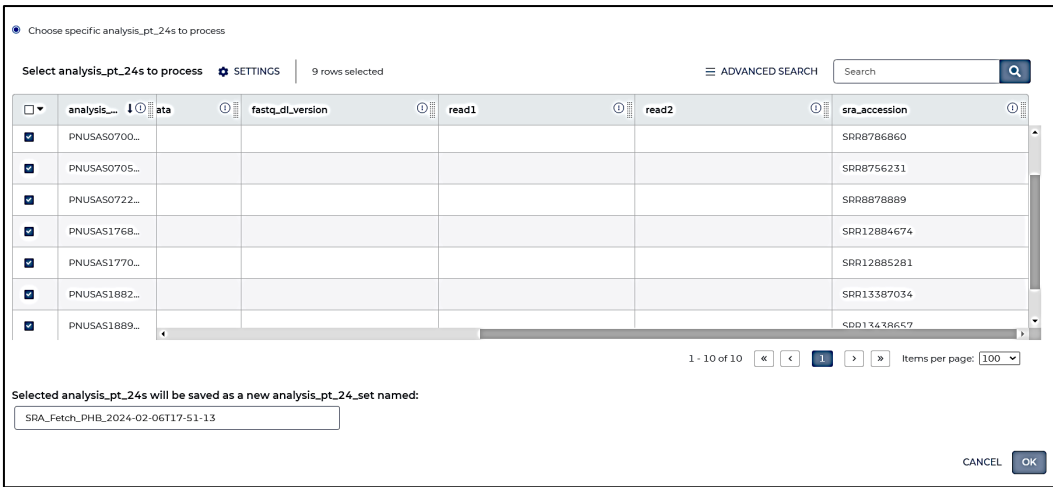


9. У розділі "Крок 1" натисніть на випадаюче меню "Select root entity type" і виберіть таблицю даних, до якої ви імпортували (кроки 1-6) tsv-файл, що містить номери приєднання до SRA для зразків, які потрібно завантажити з SRA.

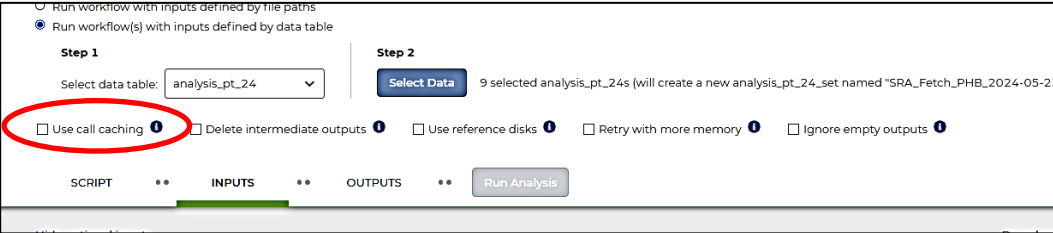


10. У розділі "Крок 2" натисніть "Вибрати дані" (скріншот вище). Ви потрапите на екран вибору вибірки.

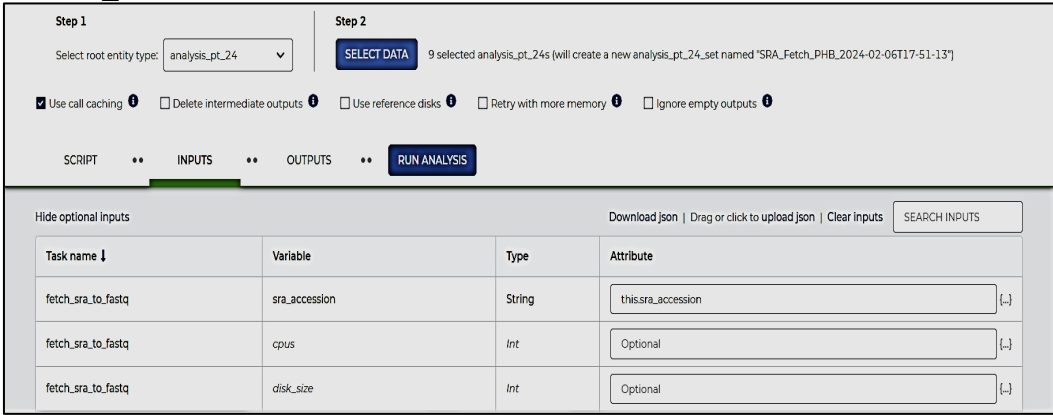
11. Встановіть прапорці навпроти зразків, які потрібно завантажити з NCBI, і натисніть "OK". Ви повернетесь на екран "SRA\_Fetch\_PHB".



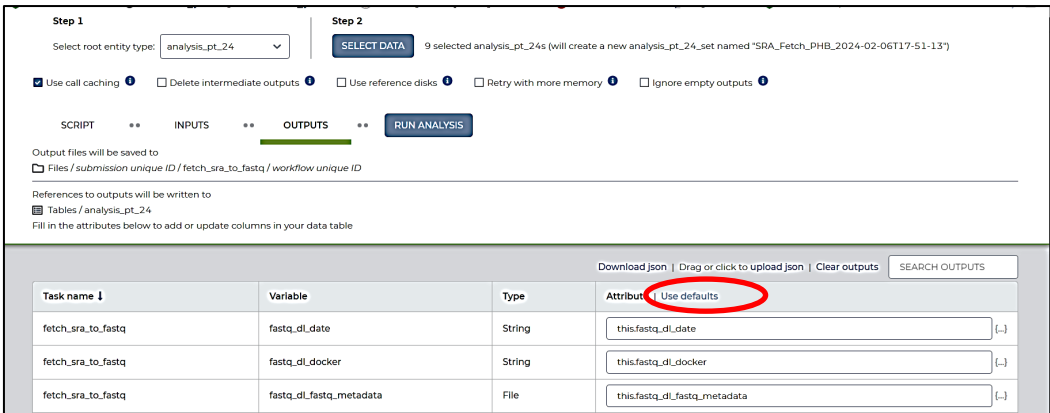
12. Зніміть прапорець з пункту " Use call catching ".



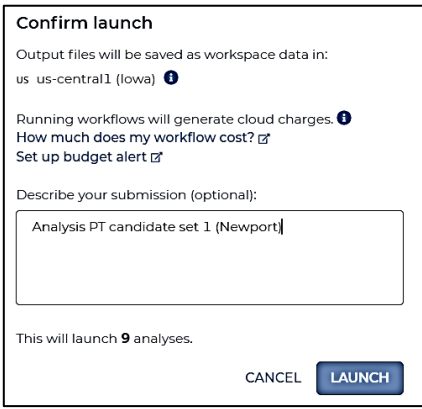
13. На вкладці "Inputs" у полі "Attribute" визначте змінну "sra\_accession": "this.sra\_accession".



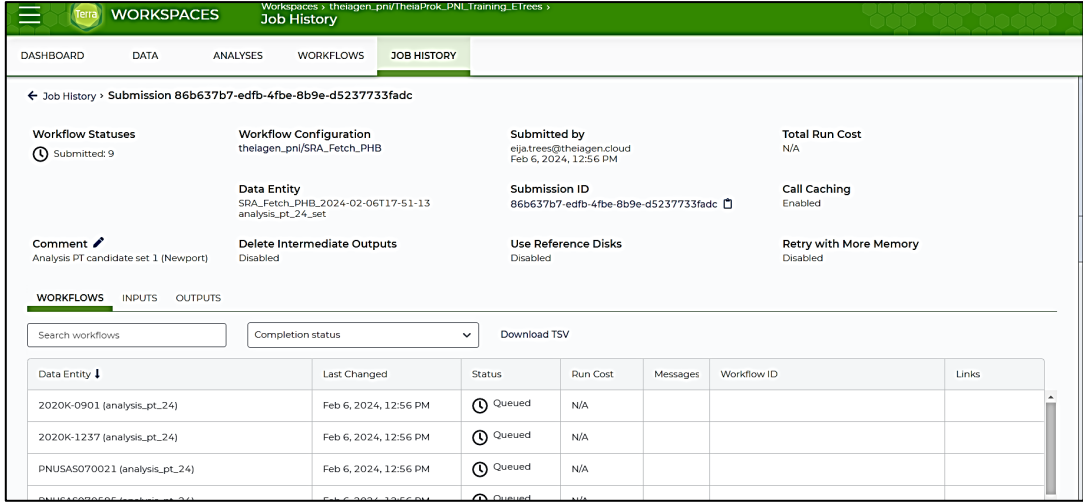
14. На вкладці "Результати" натисніть "Використовувати значення за замовчуванням", потім натисніть "Зберегти", а потім "Запустити аналіз". **ПРИМІТКА:** Кнопка "Зберегти" відображається лише в тому випадку, якщо параметри (крім ідентифікаторів зразків) змінилися з моменту попереднього подання завдання.



15. У спливаючому вікні "Підтвердити запуск" опишіть вашу заявку (необов'язково) і натисніть "Запустити".



16. Відкриється вкладка "Історія робіт", де статус поданих заявок спочатку має бути "В черзі".



17. Після завершення роботи статус буде "Виконано". На вкладці "Дані" ви побачите імена файлів FASTQ у стовпчиках "Read1" і "Read2".

**МІЖНАРОДНА СТАНДАРТНА ОПЕРАЦІЙНА ПРОЦЕДУРА PULSENET ДЛЯ АНАЛІЗУ ДАНИХ  
КОРОТКОГО ЗЧИТУВАННЯ WGS ILLUMINA З ВИКОРИСТАННЯМ ПЛАТФОРМИ TERRA.VI**

Док. No. PNID01

Вер. No. 01

Дата набуття чинності:

Сторінка 52 з 67

WORKSPACES Workspaces > thelagen\_pni/TheiaProk\_PNI\_Training\_ETrees > Job History

DASHBOARD DATA ANALYSES WORKFLOWS **JOB HISTORY**

Submission (click for details) Data entity No. of Workflows Status Submitted Submission ID Comment Actions

SRA_Fetch_PHB Submitted by <a href="mailto:elja.trees@thelagen.cloud">elja.trees@thelagen.cloud</a>	SRA_Fetch_PHB_2024-0...	9	Done	Feb 6, 2024 12:56 PM	86b637b7-edfb-47bc-8b9e-05237733fadc	Analysis PT candidate set ...	<a href="#">i</a>
BaseSpace_Fetch_PHB Submitted by <a href="mailto:curtis.kapsak@thelagen.com">curtis.kapsak@thelagen.com</a>	SRA_Fetch_PHB_2024-02-06T17:51:13 (analysis_pt_24_set)		Done	Feb 2, 2024 3:15 PM	d7a05e8e-100b-4977-ae6b5-4de5baf9a9aa	test on OC miseq run whe...	<a href="#">i</a>

WORKSPACES Workspaces > thelagen\_pni/TheiaProk\_PNI\_Training\_ETrees > Data

DASHBOARD **DATA** ANALYSES WORKFLOWS JOB HISTORY

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

TABLES

- 2300109 (2)
- CDC\_ATCC\_Sequences (42)
- CDC\_ATCC\_Sequences\_set (22)
- illumina\_PE\_v1-0-0 (25)
- analysis\_pt\_24 (10)**
- analysis\_pt\_24\_set (1)
- orange\_miseq\_specimen (20)
- orange\_miseq\_specimen\_set (1)
- pni\_training\_sample (11)

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

1 - 10 of 10 Items per page: 100

**Додаток PNID01-3: Налаштування подання таблиці даних для показників контролю якості PulseNet**

**ПРИМІТКА:** Дане налаштування можна виконати лише після того, як ви один раз запустили робочий процес *TheiaProk*.

1. На вкладці "Дані" виберіть таблицю даних, яка вас цікавить, наприклад, "CDC\_ATCC\_Sequences", а потім виберіть "Налаштування".
2. У розділі "Вибрати стовпці" слід позначити наступні метрики контролю якості:
  - 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 ⓘ

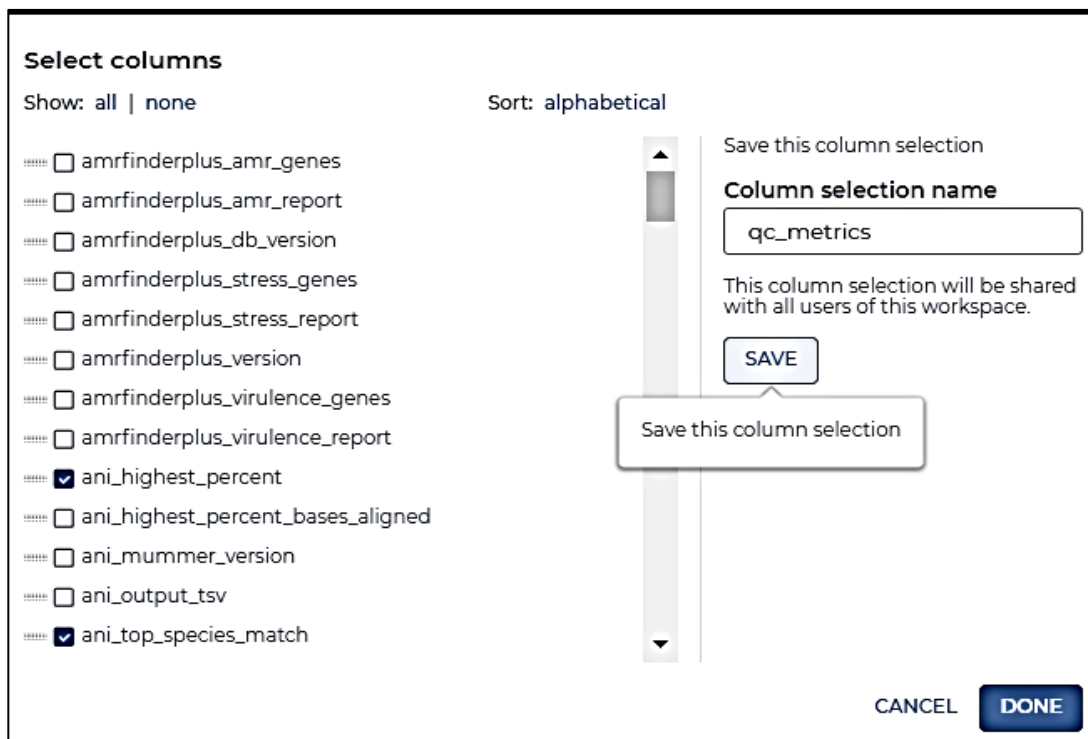
qc\_metrics ⓘ

CANCEL
DONE

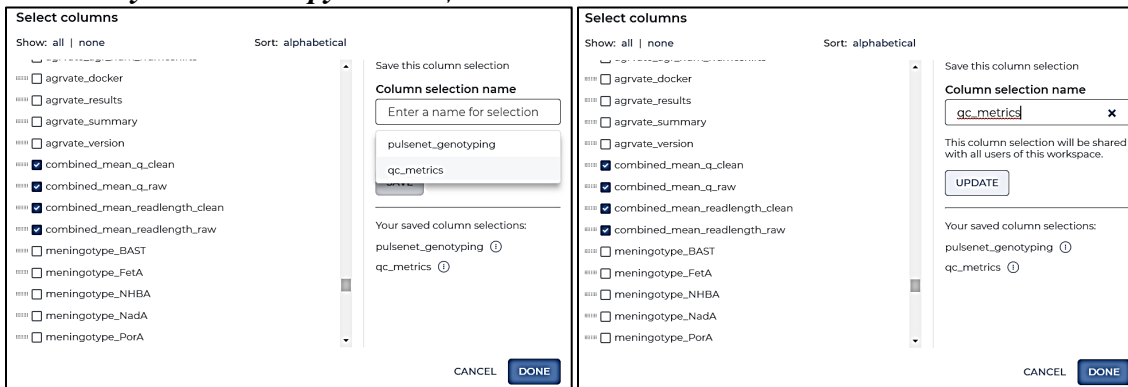
3. Натисніть "Зберегти цей вибір стовпця", назвіть вибір стовпця "qc\_metrics" і натисніть "Зберегти", а потім натисніть "Готово".

**ПРИМІТКА:** Якщо ви додаєте або видаляєте стовпці з існуючого виділення стовпців, натисніть "Зберегти це виділення стовпців", виберіть назву зі спадного меню і натисніть "Оновити".

**Створення нового виділення стовпців**



**Зміна існуючого відбору стовпців**



**Додаток PNID01-4а. Критичні показники якості PulseNet (проходження/непроходження) для подачі рутинних послідовностей**

Організм	Середнє охоплення <i>denovo</i> <sup>1</sup>	Середня якість (Q score) <sup>2</sup>	Довжина збірки (МБ)	Чисельність вторинних видів
<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>Після обрізання на основі якості (est\_coverage\_clean)

<sup>2</sup>Перед обрізанням (combined\_mean\_q\_raw)

**Додаток PNID01-4b. Етап попереднього скринінгу зчитування TheiaProk для виключення неякісних послідовностей з метою економії обчислювальних ресурсів**

Скринінгове завдання гарантує, що кількість даних про послідовності достатня для проведення геномного аналізу. Воно використовує команди bash для кількісного підрахунку зчитувань і базових пар, а також mash-схемування для оцінки розміру геному та його покриття. На кожному кроці результати оцінюються відповідно до критеріїв "пройшов/не пройшов" і порогових значень, які можуть бути визначені додатковими даними користувача.

Зразки, які не відповідають цим критеріям, не будуть далі оброблятися робочим процесом:

1. Загальна кількість зчитувань: Зразок не пройде перевірку на читання, якщо його загальна кількість прочитань менша або дорівнює min\_reads.
2. Пропорція базових пар, прочитаних у файлах прямого і зворотного зчитування: Зразок не пройде перевірку на зчитування, якщо у файлах reads1 або read2 буде менше, ніж min\_proportion базових пар.
3. Кількість базових пар: Зразок не пройде перевірку на читання, якщо кількість базових пар менша за min\_basepairs
4. Розрахунковий розмір геному: Зразок не пройде скринінг зчитування, якщо передбачуваний розмір геному менший за min\_genome\_size або більший за max\_genome\_size.
5. Розрахункове покриття геному: Зразок не пройде скринінг зчитування, якщо оцінене покриття геному менше, ніж min\_coverage.

Значення за замовчуванням:

Int min\_reads = 7472

**МІЖНАРОДНА СТАНДАРТНА ОПЕРАЦІЙНА ПРОЦЕДУРА PULSENET ДЛЯ АНАЛІЗУ ДАНИХ  
КОРОТКОГО ЗЧИТУВАННЯ WGS ILLUMINA З ВИКОРИСТАННЯМ ПЛАТФОРМИ TERRA.VI**

Док. No. PNID01

Вер. No. 01

Дата набуття чинності:

Сторінка 56 з 67

Int min\_basepairs = 2241820

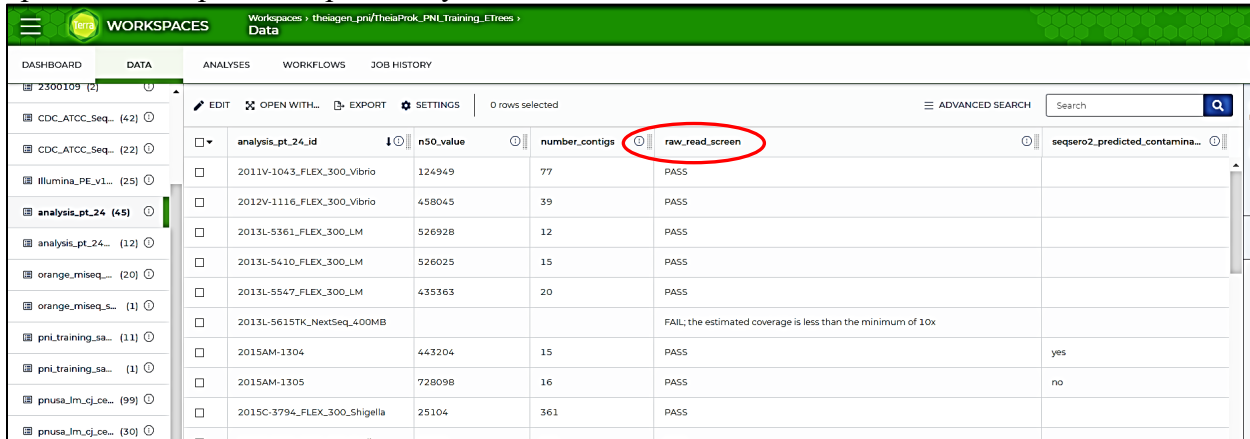
Int min\_genome\_length = 100000

Int max\_genome\_length = 18040666

Int min\_coverage = 10

Int min\_proportion = 40

У колонці "raw\_read\_screen" у розділі "QC\_metrics" буде показано, якщо зразок не пройшов попередній екран зчитування:



The screenshot shows the Terra Workspaces interface with a table of QC metrics. The table has columns for 'analysis\_pt\_24\_id', 'n50\_value', 'number\_contigs', 'raw\_read\_screen', and 'seqero2\_predicted\_contamina...'. The 'raw\_read\_screen' column is circled in red. The table contains several rows of data, including entries for 'CDC\_ATCC\_Seq...', 'illumina\_FE\_v1...', 'analysis\_pt\_24', 'orange\_miseq...', and 'pnusa\_lm\_cj.ce...'. The 'raw\_read\_screen' column shows 'PASS' for most entries, but 'FAIL; the estimated coverage is less than the minimum of 10x' for one entry.

analysis_pt_24_id	n50_value	number_contigs	raw_read_screen	seqero2_predicted_contamina...
2011V-1043_FLEX_300_Vibrio	124949	77	PASS	
2012V-1116_FLEX_300_Vibrio	458045	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	

### Додаток PNID01-5. Налаштування подання таблиці даних для аналізів генотипування PulseNet

**ПРИМІТКА:** це налаштування можна виконати лише після того, як ви один раз запустили робочий процес TheiaProk.

1. На вкладці "Дані" виберіть таблицю даних, яка вас цікавить, наприклад, "CDC\_ATCC\_Sequences", а потім виберіть "Налаштування".
2. У розділі "Вибрати стовпчики" слід позначити наступні генотипувальні аналізи:
  - 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. Натисніть "Зберегти цей вибір стовпця", назвіть вибір стовпця "pulsenet\_genotyping" і натисніть "Зберегти", а потім натисніть "Готово".

**ПРИМІТКА:** Якщо ви додаєте або видаляєте стовпці з існуючого виділення стовпців, натисніть "Зберегти це виділення стовпців", виберіть назву зі спадного меню і натисніть "Оновити".

**МІЖНАРОДНА СТАНДАРТНА ОПЕРАЦІЙНА ПРОЦЕДУРА PULSENET ДЛЯ АНАЛІЗУ ДАНИХ  
КОРОТКОГО ЗЧИТУВАННЯ WGS ILLUMINA З ВИКОРИСТАННЯМ ПЛАТФОРМИ TERRA.VI**

Док. No. PNID01

Вер. No. 01

Дата набуття чинності:

Сторінка 58 з 67

**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**

pulsenet\_genotyping

pulsenet\_genotyping

SAVE

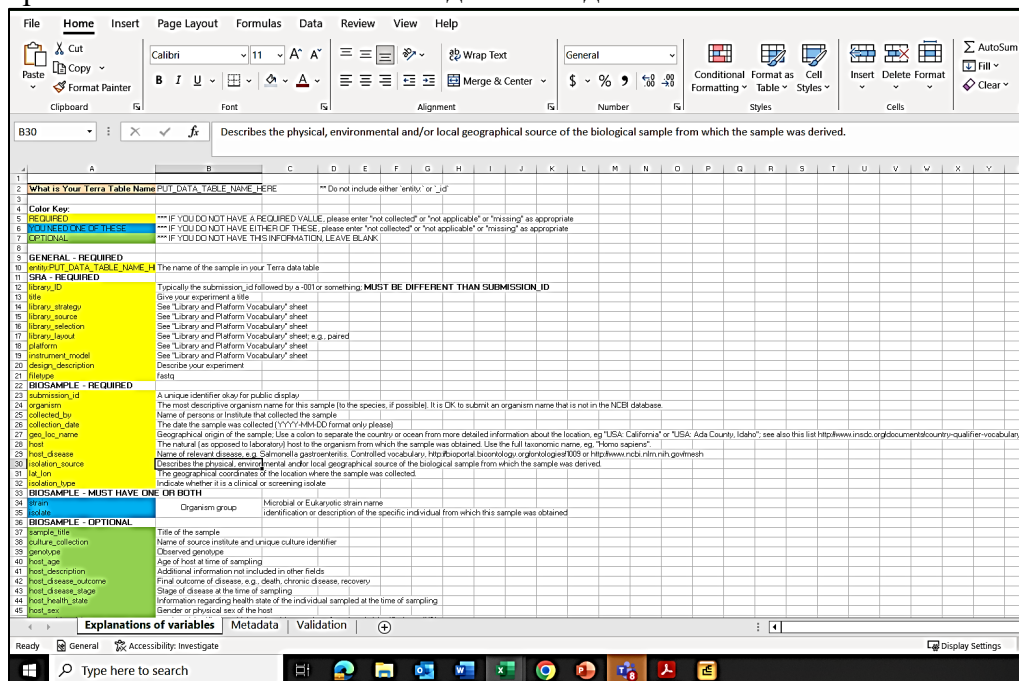
Save this column selection

qc\_metrics ⓘ

CANCEL **DONE**

### Додаток PNID01-6. Завантаження додаткових метаданих до Terra для подання до NCBI та налаштування вигляду таблиці даних для метаданих

- NCBI вимагає завантаження мінімальних метаданих для створення біозразків для послідовностей, які будуть завантажені в SRA.
- Оскільки метадані NCBI мають бути відформатовані певним чином, будь ласка, використовуйте шаблон метаданих **патогенів**, наданий компанією Theiagen, для завантаження метаданих до Terra: [https://theiagen.notion.site/Terra\\_2\\_NCBI-8f014c73acc44465a3d69cf4df93adfe](https://theiagen.notion.site/Terra_2_NCBI-8f014c73acc44465a3d69cf4df93adfe).
- Шаблон метаданих має три вкладки:
  - Перша вкладка, яка називається "Пояснення змінних", містить опис обов'язкових і необов'язкових полів.
  - Метадані для завантаження вводяться на другій вкладці "Метадані".
  - Третя вкладка "Валідація" може бути використана для перевірки правильності заповнення необхідних метаданих.



### Щоб завантажити файл метаданих патогену до Terra:

1. Заповніть обов'язкові та необов'язкові (за необхідності) поля на вкладці метаданих в електронній таблиці шаблону метаданих патогену:

**ПРИМІТКА 1:** ви можете ввести "Відсутня" для будь-якої необхідної інформації, якої у вас немає або яку ви не хочете публічно розголошувати.

**ПРИМІТКА 2:** Наведені нижче метадані є мінімальними вимогами PulseNet USA до метаданих для завантаження в NCBI, розробленими з метою захисту конфіденційності пацієнтів і цілісності поточних досліджень спалахів. З іншого



# МІЖНАРОДНА СТАНДАРТНА ОПЕРАЦІЙНА ПРОЦЕДУРА PULSENET ДЛЯ АНАЛІЗУ ДАНИХ КОРОТКОГО ЗЧИТУВАННЯ WGS ILLUMINA З ВИКОРИСТАННЯМ ПЛАТФОРМИ TERRA.BI

Док. No. PNID01

Вер. No. 01

Дата набуття чинності:

Сторінка 61 з 67

library_type	submission_id	organism	collected_by	collection_date	geo_loc_name	host	host_disease	isolation_source	lat_lon	isolation_type	strain	isolate	sample_title	culture
fastq	2017C-4936	Escherichia coli	CDC	2017-01-01	USA	Homo sapiens	Missing	Missing	Missing	clinical	2017C-4936			
fastq	2018C-4039	Escherichia coli	CDC	2018-01-01	USA	Homo sapiens	Missing	Missing	Missing	clinical	2018C-4039			
fastq	2019C-3204	Shigella sonnei	CDC	2018-01-01	USA	Homo sapiens	Missing	Missing	Missing	clinical	2019C-3204			

## Клінічні (людського походження) та неклінічні послідовності

entity_quality_control_id	library_id	title	library_strategy	library_source	library_selection	library_layout	platform	instrument_model	design_description	filetype	submission_id	organism
2019C-3238	2019C-3238-001	PulseNet WGS	GENOMIC	RANDOM	paired	ILLUMINA Illumina HiSeq 2500	Missing	fastq	2019C-3238	Escherichia coli		
2015C-3887	2015C-3887-001	PulseNet WGS	GENOMIC	RANDOM	paired	ILLUMINA Illumina HiSeq 2500	Missing	fastq	2015C-3887	Escherichia coli		
2017C-4938	2017C-4938-001	PulseNet WGS	GENOMIC	RANDOM	paired	ILLUMINA Illumina HiSeq 2500	Missing	fastq	2017C-4938	Escherichia coli		
2013L-5357-LRM4update	2013L-5357-LRM4update-001	PulseNet WGS	GENOMIC	RANDOM	paired	ILLUMINA Illumina MiSeq	Missing	fasta	2013L-5357-LRM4update	Listeria monocytogenes		

organism	collected_by	collection_date	geo_loc_name	host	host_disease	isolation_source	lat_lon	isolation_type	strain	isolate
Escherichia coli	CDC	2017-01-01	USA	Missing	Missing	Missing	Missing	clinical	2017C-4936	
Escherichia coli	CDC	2018-01-01	USA	Missing	Missing	Missing	Missing	clinical	2018C-4039	
Escherichia coli	CDC	2018-01-01	USA	Missing	Missing	Missing	Missing	clinical	2019C-3204	
Listeria monocytogenes	CDC	2013-07-01	USA:WI	Missing	Missing	cheese	Missing	food	2013L-5357-LRM4update	

2. Перейдіть на вкладку "Валідація", щоб підтвердити правильність заповнення шаблону метаданих.

**This tab is intended to help ensure your metadata fields are valid.**  
Not all issues can be captured here; please carefully examine the Explanations of variables sheet for detailed information.  
Please be aware that even if your metadata passes this validation check, it does not ensure 100% valid metadata fields.

Did you correctly fill in your Terra table name on the Explanation of Variables tab?		No, please update this field.			
CHECKING REQUIRED FIELDS...					
Row #	Are the library_id and submission_id different?	Do all required fields have values?	Does geo_loc_name have the proper format?	Is the collection_date valid?	Does either the strain or isolate fields have values?
8	2 Yes	Yes	No	Yes	Yes
9	3 Yes	Yes	No	Yes	Yes
10	4 Yes	Yes	No	Yes	Yes
11					
12					

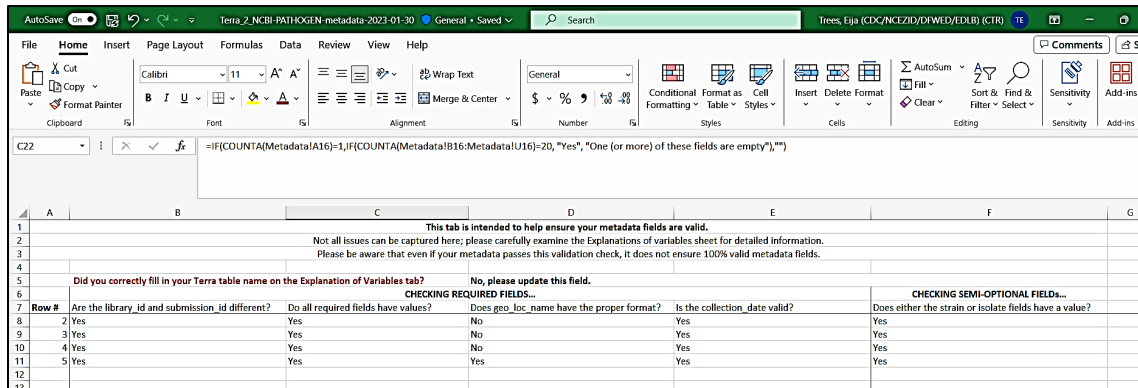
**МІЖНАРОДНА СТАНДАРТНА ОПЕРАЦІЙНА ПРОЦЕДУРА PULSENET ДЛЯ АНАЛІЗУ ДАНИХ  
КОРОТКОГО ЗЧИТУВАННЯ WGS ILLUMINA З ВИКОРИСТАННЯМ ПЛАТФОРМИ TERRA.BI**

Док. No. PNID01

Вер. No. 01

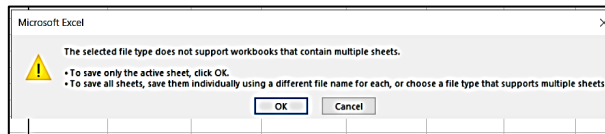
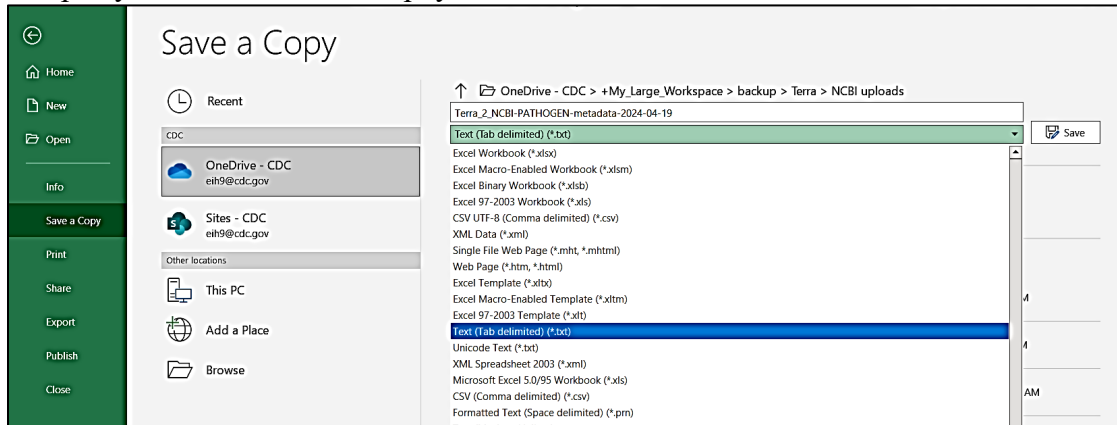
Дата набуття чинності:

Сторінка 62 з 67



**ПРИМІТКА:** валідація належного формату географічного розташування пройде тільки в тому випадку, якщо вказана країна-джерело і більш детальне місцезнаходження. Таким чином, валідація не пройде для клінічних ізолятів, для яких вказана лише країна походження. Однак NCBI прийме країну походження як єдине географічне місцезнаходження.

- Збережіть заповнений шаблон метаданих у форматі tsv (з табуляцією). Натисніть "ОК" у спливаючому вікні, в якому зазначено, що вибраний тип файлу не підтримує книги з кількома аркушами.



- У вкладці "Дані" робочого простору Terra натисніть "Імпортувати дані" і виберіть "Завантажити TSV" у випадяючому меню.

quality_control_id	read1	read2
2017C-4938	14_L001_R1_001.fastq.gz	2017C-4938-D00290
2017C-4938	02_L001_R1_001.fastq.gz	2017C-4938-D00290
2018C-4039	94_L001_R1_001.fastq.gz	2018C-4039-D00290
2018C-4709-L-A1-USA-CDC-pcl	L001_R1_001.fastq.gz	2018C-4709-L-A1-M
2018C-4709-L-A2-USA-CDC-pcl	L001_R1_001.fastq.gz	2018C-4709-L-A2-M

5. У спливаючому вікні "Імпорт даних таблиці" на вкладці "Імпорт файлів" клацніть посередині, щоб вибрати файл tsv.

**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-central1 (Iowa)

**FILE IMPORT**    TEXT IMPORT

Select the TSV file containing your data:

Drag or Click to select a .tsv file

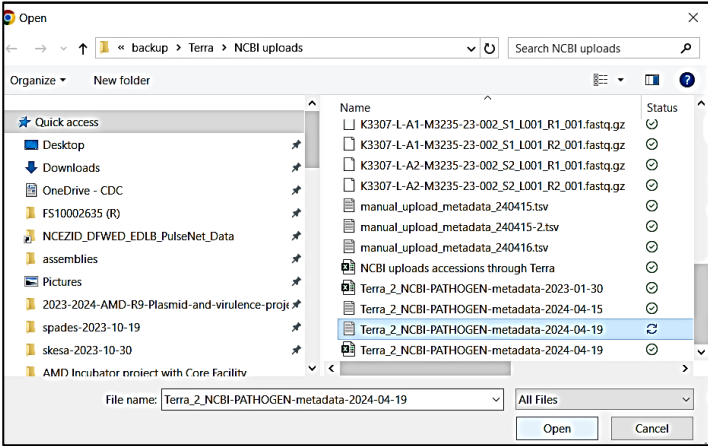
Selected File: None

**TSV file templates**

Download sample\_template.tsv

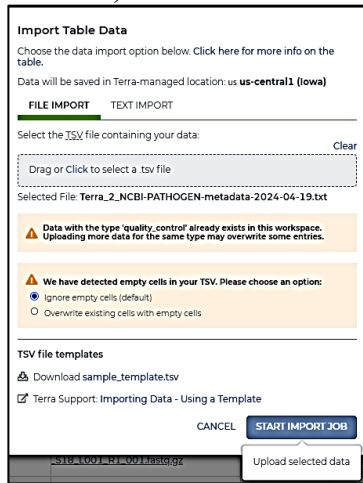
Terra Support: Importing Data - Using a Template

6. Перейдіть до місця, де зберігається файл метаданих tsv, виберіть його і натисніть "Відкрити".



7. У спливаючому вікні "Імпортувати дані таблиці" ви побачите попередження про те, що дані вже існують у відповідній таблиці даних, і завантаження нових даних до неї може призвести до перезапису наявних даних. Також з'явиться попередження,

якщо файл метаданих tsv не містить однакової інформації (деякі дані відсутні для деяких послідовностей) для всіх записів. Натисніть на кнопку "Почати імпорт".



8. Після завершення завантаження ви побачите бажані поля метаданих, заповнені в таблиці даних.

quality_control_id	collected_by	collection_date	design_description	filetype
2017C-3830-LRM4update				
2017C-4936	CDC	2017-01-01	Missing	fastq
2017C-4938				
2018C-4039	CDC	2018-01-01	Missing	fastq
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	CDC	2018-01-01	Missing	fastq
2019C-3238				
92-01-L-A1	CDC	Missing	Missing	fastq
92-01-L-A2				

### Створіть окреме подання метаданих для таблиці даних

1. На вкладці "Дані" виберіть таблицю даних, яка вас цікавить, а потім виберіть "Налаштування".
2. У розділі "Вибрати стовпці" позначте всі потрібні стовпці метаданих.
  - a. Рекомендовані метадані для подачі в NCBI:
    - 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. Додаткова корисна інформація про послідовність
  - i. read1 (ім'я файлу R1 FASTQ)
  - ii. read2 (ім'я файлу R2 FASTQ)
  - iii. assembly\_fasta (місце розташування збірки, згенерованої Terra)
  - iv. Якщо ви завантажуєте дані безпосередньо з 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

Select columns

Show: all | none      Sort: alphabetical

design\_description

filetype

geo\_loc\_name

host

host\_disease

instrument\_model

isolation\_source

isolation\_type

lat\_lon

library\_ID

library\_layout

library\_selection

library\_source

library\_strategy

SAVE THIS COLUMN SELECTION

Your saved column selections:

Metadata ⓘ

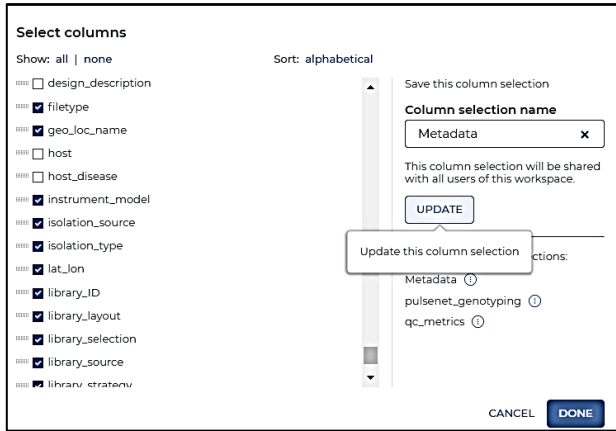
pulsenet\_genotyping ⓘ

qc\_metrics ⓘ

CANCEL    DONE

- 3. Натисніть "Зберегти цей вибір стовпця".
- 4. Назвіть вибір стовпця "Метадані" і натисніть "Зберегти" та "Готово".

**ПРИМІТКА:** Якщо ви додаєте або видаляєте стовпці з існуючого виділення стовпців, натисніть "Зберегти це виділення стовпців", виберіть назву зі спадного меню і натисніть "Оновити".



5. Тепер у таблиці даних мають бути видимими лише потрібні стовпці метаданих

The screenshot shows the Terra Data table. The table has columns: quality\_control\_id, sample\_accession, collected\_by, collection\_date, filetype, and geo. The data rows are as follows:

quality_control_id	sample_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					

The table is part of a web interface with a green header 'WORKSPACES' and a navigation bar with 'DASHBOARD', 'DATA', 'ANALYSES', 'WORKFLOWS', and 'JOB HISTORY'. The 'DATA' tab is active. The table has 94 rows in total, with 100 items per page.

### Додаток PNID01-7: Огляд робочого процесу TheiaProk для визначення бактеріальних характеристик

