# BTEP course

Center for Cancer Research

**BTEP**

Bioinformatics Training
& Education Program

Alexandra L Emmons Ph.D. & Joe Wu Ph.D.
BTEP/GAU/CCR/NCI/NIH - email ncibtep@mail.nih.gov
Bioinformatics Training and Education Program

# Table of Contents

# Course Overview

Partek Flow is a start-to-finish solution for analyzing high dimensional multi-omics sequencing data. It is a point-and-click software and is suitable for those who wish to avoid the steep learning curve associated with analyzing sequencing data through command line and/or code. At NIH, Partek Flow is hosted on the Biowulf *(https://partekflow.cit.nih.gov)* high performance computing cluster (HPC). Researchers interact with the software through a web browser using a URL supplied by Biowulf once a Biowulf and Partek Flow account has been set up. This enables investigators to take advantage of the compute power offered by HPC while using a graphical user interface to construct a sequencing data analysis workflow. Partek Flow enables the creation of publication quality visualizations.

This in-person and hands-on training will introduce participants to single cell RNA sequencing (move this to beginning) analysis on Partek Flow using a single human PBMC sample count matrix. Participants will also learn how to access and methods for transferring data to the Partek Flow server at NIH. Skills acquired from this class are applicable to analyzing other types of sequencing data in Partek Flow.

Date: March 5, 2024 Time: 2 - 4 PM Location: Building 35A Room 620/630

## Learning Objectives

After this class, participants will

- Become familiar analyzing single cell RNA sequencing data using Partek Flow including
    - Importing of data into a Partek Flow project
    - QA/QC, filtering, and normalizing of single cell RNA data
    - Performing cell type classification based on gene expression
    - Performing differential expression and pathway analysis
    - Producing visualization (PCA, UMAP, tSNE, dotplot, volcano plot, hierarchical clustering etc.)
- Know how to access Partek Flow
- Know how to sign onto the NIH Partek Flow server
- Be able to transfer data from NCI CCR Sequencing Facility Data Management Environment to their Biowulf Partek Flow folder

## Link to documents from Partek

Class documents provided by Partek *(https://bioinformatics.ccr.cancer.gov/btep/wp-content/uploads/sites/2/PartekFlowSCTrainingHandout_MAR_2024.pdf)*

# Accessing Partek Flow at NIH and tips for data transfer

## Learning objectives

After consulting this guide, participants will

- Know how to access Partek Flow at NIH.
- Be able to transfer data from NCI CCR Sequencing Facility Data Management Environment to their Biowulf Partek Flow folder.

## Instructions for accessing Partek Flow

NCI researchers can find instructions for accessing Partek Flow at https:// bioinformatics.ccr.cancer.gov/btep/partek-flow-bulk-and-single-cell-rna-seq-data-analysis/ *(https://bioinformatics.ccr.cancer.gov/btep/partek-flow-bulk-and-single-cell-rna-seq-data-analysis/)*. But the things needed are

- A Biowulf (The High Performance Computing cluster) account — see here for information about how to obtain a HPC account *(https://hpc.nih.gov/docs/accounts.html)*.
- A /data directory on Biowulf with enough disk space to hold their Partek Flow files — please fill out this online form *(https://hpc.nih.gov/dashboard/storage_request.php)* if you do not already have a /data directory or if you require more disk space.
- A Partek Flow account created for them — please contact staff@hpc.nih.gov.

Once these steps have been accomplished, Partek Flow is available at https:// partekflow.cit.nih.gov/flow *(https://partekflow.cit.nih.gov/flow)*.

## The Partek Flow folder on Biowulf

HPC staff will create a folder called "PartekFlow" in the user's Biowulf data directory. This folder will hold all Partek Flow projects.

## Transferring data from NCI CCR Sequencing Facility to Partek Flow on Biowulf

Those researchers who used the NCI CCR Sequencing Facility *(https:// bioinformatics.ccr.cancer.gov/docs/resources-for-bioinformatics/raw_data_from_cores/)* to get

sequencing done will receive a link to their data. This data can be transferred to the "PartekFlow" folder on Biowulf using Globus. The steps for setting up a Globus endpoint for the Biowulf "PartekFlow" folder can be found at https://partekflow.cit.nih.gov/#upload_globus *(https://partekflow.cit.nih.gov/#upload_globus)*. The embedded PDF shows how to connect the sequencing facility's data management environment to a Globus endpoint.

For those who have not setup a Globus account, refer to https://hpc.nih.gov/docs/globus/setup.php *(https://hpc.nih.gov/docs/globus/setup.php)* for instructions.

> **Tip**
>
> If following the Biowulf instructions for creating a Globus endpoint for the "PartekFlow" folder, it will be a good idea to use subdirectories for data generated for different experiments. This exercise will use a subdirectory called fnl_example_single_cell_fastq.

## Sign onto Globus

Information regarding Globus and how to obtain it can be found at https://hpc.nih.gov/docs/globus/setup.php *(https://hpc.nih.gov/docs/globus/setup.php)*.

For those Globus already setup, goto https://www.globus.org *(https://www.globus.org)* to log in by clicking on the "LOG IN" icon at the top right of the pages.

After clicking on the log in button, select organziational affiliation, which is National Institutes of Health in this example.



Click on "Continue" when the organizational affiliation has been selected.

After clicking "Continue", users will be brought the to Globus interface where file transfers are managed. Clicking on "COLLECTIONS" and then check "ADMINISTER BY YOU" will reveal several endpoints including on that points to that for the instructor's Biowulf "PartekFlow" folder labeled "example fastq from fnl sf dme to biowulf partek flow" to see the overview. Note the UUID.
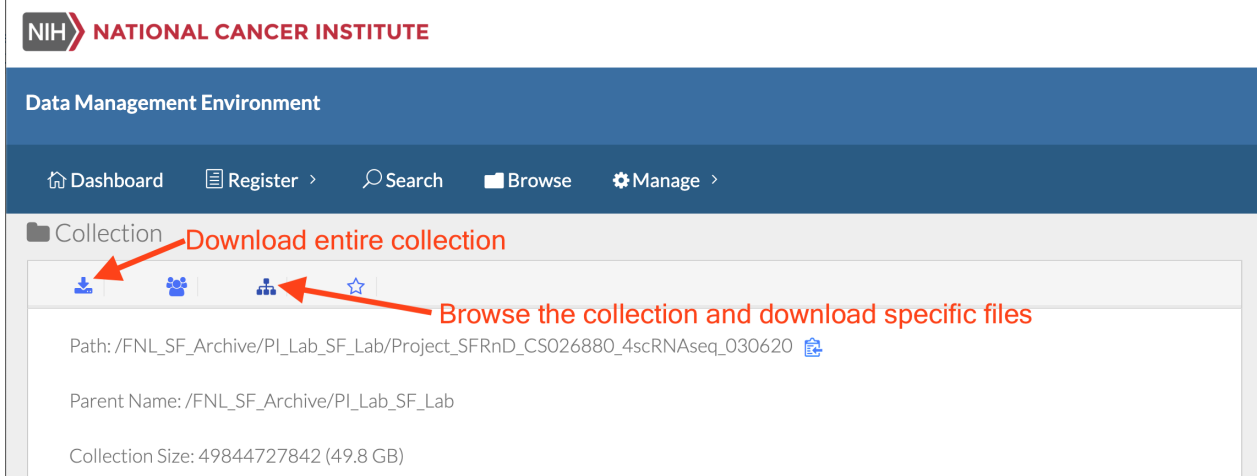


## NCI CCR Sequencing Facility Data Management Environment

The NCI CCR Sequencing Facility will send researchers a link to the data, which is stored in their Data Management Environment (DME). Again, instructions for connecting DME to a Globus endpoint on Biowulf are in the embedded PDF.

Users are able to download an entire collection of data or browse the collection and download a subset.

**NIH⟩ NATIONAL CANCER INSTITUTE**

**Data Management Environment**

⌂ Dashboard    🗐 Register ›    🔍 Search    ■ Browse    ⚙ Manage ›

📁 Collection   Download entire collection

⬇    👥    🔗    ☆

Browse the collection and download specific files

Path: /FNL_SF_Archive/PI_Lab_SF_Lab/Project_SFRnD_CS026880_4scRNAseq_030620 📑

Parent Name: /FNL_SF_Archive/PI_Lab_SF_Lab

Collection Size: 49844727842 (49.8 GB)

This example will browse the collection and download the N_1395BL_NextGEM_count.tar.

At the subsequent page, enter the UUID for the Biowulf Partek Flow Globus endpoint and "/" for the path. Then click "Download".

**NIH⟩ NATIONAL CANCER INSTITUTE**

**Data Management Environment**

⌂ Dashboard    🗐 Register ›    🔍 Search    ■ Browse    ⚙ Manage ›

⬇ Download

This page allows you to download the selected collection to a Globus endpoint, an AWS S3 bucket, Google Drive, Google Cloud or dbGaP.

Selected Collection:

/FNL_SF_Archive/PI_Lab_SF_Lab/Project_SFRnD_CS026880_4scRNAseq_030620/Flowcell_HYH2JBGXC/Sample_N_1395BL_NextGEM

◉ Globus   ○ AWS S3   ○ Google Drive   ○ Google Cloud   ○ dbGaP

To download to Globus, DME should be provided write access to the targeted Globus endpoint. Refer to Preparing to Use Globus with DME for instructions.

Globus Endpoint UUID: Obtain from Globus (Optional)

> 9a06b185-9d72-492f-a568-5e958ffe3a9f

Globus Endpoint (Destination) Path:

> /

Download

**NIH》 NATIONAL CANCER INSTITUTE**

**Data Management Environment**

⌂ Dashboard    ☰ Register ›    🔍 Search    ■ Browse    ⚙ Manage ›

⬇ Download

This page allows you to download the selected collection to a Globus endpoint, an AWS S3 bucket, Google Drive, Google Cloud or dbGaP.

Selected Collection:

/FNL_SF_Archive/PI_Lab_SF_Lab/Project_SFRnD_CS026880_4scRNAseq_030620/Flowcell_HYH2JBGXC/Sample_N_1395BL_NextGEM

Asynchronous download request is submitted successfully! Task Id: 05ec5175-6216-44cb-861e-887902e83485

⦿ Globus    ◯ AWS S3    ◯ Google Drive    ◯ Google Cloud    ◯ dbGaP

To download to Globus, DME should be provided write access to the targeted Globus endpoint. Refer to Preparing to Use Globus with DME for instructions.

Globus Endpoint UUID: Obtain from Globus (Optional)

9a06b185-9d72-492f-a568-5e958ffe3a9f

Globus Endpoint (Destination) Path:

/

---

**NIH》 NATIONAL CANCER INSTITUTE**

**Data Management Environment**

⌂ Dashboard    ☰ Register ›    🔍 Search    ■ Browse    ⚙ Manage ›

|                    |
|--------------------|
| Notifications      |
| Download Tasks     |
| Registration Tasks |

⬇ Download

This page allows you to download the selected collection to a Globus endpoint, an AWS S3 bucket, Google Drive, Google Cloud or dbGaP.

Selected Collection:

/FNL_SF_Archive/PI_Lab_SF_Lab/Project_SFRnD_CS026880_4scRNAseq_030620/Flowcell_HYH2JBGXC/Sample_N_1395BL_NextGEM

Asynchronous download request is submitted successfully! Task Id: 05ec5175-6216-44cb-861e-887902e83485

⦿ Globus    ◯ AWS S3    ◯ Google Drive    ◯ Google Cloud    ◯ dbGaP

To download to Globus, DME should be provided write access to the targeted Globus endpoint. Refer to Preparing to Use Globus with DME for instructions.

Globus Endpoint UUID: Obtain from Globus (Optional)

9a06b185-9d72-492f-a568-5e958ffe3a9f

Globus Endpoint (Destination) Path:

/

Globus will also send an email to the user's NIH email account after transfer has been completed.

These files will show up on Biowulf as well. These files will need to be unpacked using `tar -xvf`.

```
[wuz8@biowulf fnl_example_single_cell_fastq]$ ls -1 *.tar
N_1395BL_NextGEM_count.tar
```

# Importing data to Partek Flow project

Log into Partek Flow at https://partekflow.cit.nih.gov/flow *(https://partekflow.cit.nih.gov/flow)*. This example will use the nci_ccr_sf_example_scrna project, so click on it.



Click on "Add data".

Select Single cell, scRNA-Seq, and check 10x Genomics Cell Ranger counts h5. Then click Next.



Navigate PartekFlow, globus, fnl_example_single_cell_fastq, N_1395BL_NextGEM, outs and select the filtered_feature_bc_matrix.h5 file and then click Next at the bottom of the screen.



In the subsequent page, provide an informative sample name and select the appropriate assembly. Then click Finish.

Flow 🪵 Queue ⌄ Projects ⌄ Help ⌄

**Home > nci_ccr_sf_example_scrna > Initial import > Import sample files > File format options** ▶

**Sample names**

| ☑ | Sample name | Files | Cells | Features |
|---|---|---|---|---|
| ☑ | N_1395BL_NextGEM | filtered_feature_bc_matrix.h5 | 7824 | 36601 |

**Feature annotation**

☑ **Use annotation file**
Select the file that has been used to generate the feature counts (e.g. gene or protein information).

**Assembly**

Homo sapiens (human) - hg38 ⌄

**Annotation model**

Ensembl Transcripts release 110 (jstoddard) ⌄

**Primary feature identifier**

◉ Feature name (Values: MIR1302-2HG, FAM138A, OR4F5, AL627309.1, AL627309....)

○ Feature ID (Values: ENSG00000243485, ENSG00000237613, ENSG00000186092,...)

**Deduplication method**
If the feature ID is not unique, the feature will be summarized by the selected method.

◉ Mean ○ Maximum ○ Sum

**Count value format**

◉ Raw count   ○ Normalized count with log base  None ⌄

**Report**

◉ All features ○ Features with non-zero values across all samples

☑ **Cells with total read count at least**
A low total read count threshold will result in a large number of cells which might take a long time to import

400 ⌃⌄

[Back] [Finish]

When import is done there will be a "Single cell counts" data node in the Analyses window.