# GDC scRNA-Seq Analysis

**19 May 2025**

Bill Wysocki, Ph.D – Director of User Services
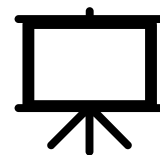Zhenyu Zhang, Ph.D – Director of Bioinformatics
  Center for Translational Data Science, University of Chicago

Xin Zhou, Ph.D - Director of Data Visualization
  Computational Biology Department, St. Jude Children's Research Hospital

**NIH** NATIONAL CANCER INSTITUTE

# Webinar Logistics

- *Webinar will be recorded*

- *Recording and slides will be made available soon*

- *Type any questions in the Q&A panel – they will be addressed at the end*
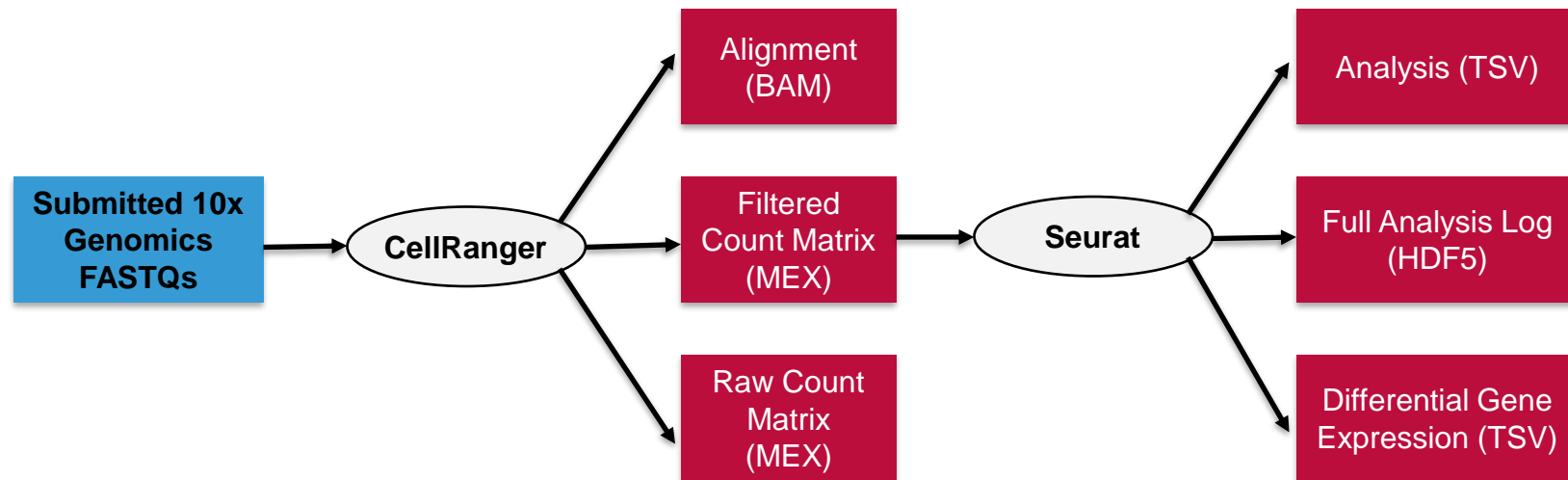
# Agenda

1. *Overview of GDC scRNA-Seq harmonization workflows*

2. *GDC scRNA-Seq file overview*

3. *GDC Single Cell RNA-Seq tool*

4. *GDC scRNA-Seq expression API*

5. *Questions*

# Overview of GDC scRNA-Seq Harmonization Workflows

# scRNA-Seq Harmonization in the GDC



- All released scRNA-Seq files are open-access, except for the BAMs

# scRNA-Seq File Overview

NIH | NATIONAL CANCER INSTITUTE

## GDC Single Cell RNA-Seq Support

**Date:** Monday, September 27, 2021
**Time:** 2:00 PM - 3:00 PM EDT
**Location:** Web Conference (See WebEx information below)

Single-cell RNA-Seq (scRNA-Seq) is a valuable tool for studying tumor heterogeneity and also the microenvironment, which may be composed of a distinct cancer subclones along with immune and other cell types.

Read More

# scRNA-Seq File Formats – Gene Expression

File.tar.gz - MEX Format

- **barcodes.tsv.gz** - a list of barcodes for this sample

- **features.tsv.gz** - a list of each gene id and name

- **matrix.mtx.gz** - the gene expression matrix

- Note: If you are using a Mac, unzipping the three internal files may not work in the browser.  Use the following command in terminal: `gunzip barcodes.tsv.gz`

```
ENSG00000243485.3      RP11-34P13.3    Gene Expression
ENSG00000274890.1      MIR1302-9       Gene Expression
ENSG00000237613.2      FAM138A Gene Expression
ENSG00000268020.3      OR4G4P  Gene Expression
ENSG00000240361.1      OR4G11P Gene Expression
ENSG00000186092.4      OR4F5   Gene Expression
ENSG00000238009.5      RP11-34P13.7    Gene Expression
ENSG00000239945.1      RP11-34P13.8    Gene Expression
ENSG00000233750.3      CICP27  Gene Expression
ENSG00000268903.1      RP11-34P13.15   Gene Expression
ENSG00000269981.1      RP11-34P13.16   Gene Expression
ENSG00000239906.1      RP11-34P13.14   Gene Expression
ENSG00000241860.5      RP11-34P13.13   Gene Expression
ENSG00000222623.1      RNU6-1100P      Gene Expression
ENSG00000241599.1      RP11-34P13.9    Gene Expression
ENSG00000279928.1      F0538757.3      Gene Expression
ENSG00000279457.2      F0538757.2      Gene Expression
ENSG00000273874.1      MIR6859-2       Gene Expression
ENSG00000275135.1      F0538757.1      Gene Expression
ENSG00000228463.7      AP006222.2      Gene Expression
ENSG00000241670.3      AP006222.1      Gene Expression
ENSG00000236679.2      RP4-669L17.1    Gene Expression
ENSG00000236743.1      RP5-857K21.15   Gene Expression
ENSG00000236601.1      RP4-669L17.2    Gene Expression
ENSG00000237094.10     RP4-669L17.10   Gene Expression
ENSG00000269732.1      WBP1LP7 Gene Expression
ENSG00000278566.1      OR4F29  Gene Expression
ENSG00000224813.2      RP4-669L17.4    Gene Expression
ENSG00000233653.3      CICP7   Gene Expression
ENSG00000250575.1      RP4-669L17.8    Gene Expression
ENSG00000278757.1      U6      Gene Expression
ENSG00000231709.1      RP5-857K21.1    Gene Expression
ENSG00000235146.2      RP5-857K21.2    Gene Expression
ENSG00000239664.2      RP5-857K21.3    Gene Expression
ENSG00000230021.6      RP5-857K21.4    Gene Expression
ENSG00000223659.1      RP5-857K21.5    Gene Expression
ENSG00000225972.1      MTND1P23        Gene Expression
ENSG00000225630.1      MTND2P28        Gene Expression
ENSG00000276171.1      AC114498.1      Gene Expression
ENSG00000237973.1      RP5-857K21.6    Gene Expression
ENSG00000278791.1      MIR6723 Gene Expression
ENSG00000229344.1      RP5-857K21.7    Gene Expression
ENSG00000240409.1      MTATP8P1        Gene Expression
ENSG00000248527.1      MTATP6P1        Gene Expression
ENSG00000198744.5      RP5-857K21.11   Gene Expression
ENSG00000268663.1      WBP1LP6 Gene Expression
ENSG00000273547.1      OR4F16  Gene Expression
```

# scRNA-Seq File Formats – Analysis

> TSV Format – Columns:

- Contains information on clustering and visualization.

- cell_barcode
- read_count
- gene_count
- seurat_cluster
- Coordinates for UMAP, tSNE and PCA

# scRNA-Seq File Formats – Differential Gene Expression

TSV Format – Columns:

- Contains information on the DGE between clusters for differentially expressed genes.

- cluster
- gene
- avg_logFC – average in fold-change
- p_val – p-value
- p_val_adj – adjusted p-value
- prop_in_cluster – proportion of cells in the cluster
- prop_out_cluster

# scRNA-Seq File Formats – Full Seurat Analysis Log

HDF5 Format – Loom File:

- Contains the full analysis from Seurat
- Stores on-disk rather than in-memory.
- Can be processed with several R or Python packages (loomR or loompy)

# *Tool Demo*

# *API Demo*

NIH NATIONAL CANCER INSTITUTE

# *Questions?*

U.S. Department of Health & Human Services

National Institutes of Health | National Cancer Institute

https://www.cancer.gov/

1-8 0 0 - 4 - C A N C E R          Produced May 2025