# Introduction to scATAC

## Submitted To:

Dr. Riasat Azim
Assistant Professor, Dept. of CSE

## Submitted By:

Israt Jahan Khan (ID: 011201142)
Md. Fahim Bin Amin (ID: 011201158)

Date: 08 October 2024

# Introduction to scATAC

scATAC stands for single-cell Assay for Transposase-Accessible Chromatin and is a genomic technique used to analyze chromatin accessibility at the single-cell level. It is an adaptation of the bulk ATAC-seq method, which allows researchers to study the open regions of chromatin across an entire cell population.

Key Features and Aspects

- Single-cell resolution: Unlike bulk ATAC-seq, scATAC provides information about chromatin accessibility in individual cells, allowing for the identification of cell-type-specific regulatory elements and heterogeneity within cell populations.

- Chromatin accessibility: The technique measures which regions of the genome are "open" or accessible to regulatory proteins, such as transcription factors. These open regions are often associated with active gene regulation.

- Transposase-based method: scATAC uses a hyperactive Tn5 transposase to insert sequencing adapters into accessible regions of the chromatin.

- Sparse data: Due to the limited amount of DNA in a single cell, scATAC data is typically sparse, with many regions having zero reads in individual cells.

- Combinatorial indexing: Many scATAC protocols use combinatorial indexing to process thousands of cells in parallel, increasing throughput.

- Integration with other single-cell methods: scATAC data can be integrated with single-cell RNA sequencing (scRNA-seq) data to provide a more comprehensive view of cellular states and gene regulation.

- Applications: scATAC is used to study cellular heterogeneity, identify cell types and states, map regulatory elements, and understand gene regulation dynamics in complex tissues or during developmental processes.

- Bioinformatics challenges: Analyzing scATAC data requires specialized computational methods to handle the sparsity and high dimensionality of the data.

scATAC has become an important tool in genomics and epigenomics research, providing insights into gene regulation at unprecedented resolution and scale. It is particularly useful in fields such as developmental biology, cancer research, and immunology, where understanding cellular heterogeneity is crucial.

# Popular Datasets for Experimental Purposes

There are several popular datasets available for experimenting with scATAC-seq analysis. These datasets are often used for benchmarking, method development, and learning purposes. Here are some well-known datasets:

1. 10x Genomics Datasets:

10x Genomics provides several publicly available scATAC-seq datasets, including:
- PBMC (Peripheral Blood Mononuclear Cells)
- Mouse Brain Cells
- Human Brain Cells
- Available at: https://support.10xgenomics.com/single-cell-atac/datasets

2. Buenrostro et al. (2018) Dataset:

- Contains scATAC-seq data from human hematopoiesis
- Published in Nature: "Integrated Single-Cell Analysis Maps the Continuous Regulatory Landscape of Human Hematopoietic Differentiation"
- Available through GEO: GSE96772 (https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE96772)

3. Cusanovich et al. (2018) Dataset:

- Large-scale single-cell chromatin accessibility profiles from mouse tissues
- Published in Cell: "A Single-Cell Atlas of In Vivo Mammalian Chromatin Accessibility"
- Available through GEO: GSE111586 (https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE111586)

4. Satpathy et al. (2019) Dataset:

- scATAC-seq data from human immune cells
- Published in Nature Biotechnology: "Massively parallel single-cell chromatin landscapes of human immune cell development and intratumoral T cell exhaustion"
- Available through GEO: GSE129785 (https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE129785)

5. Lareau et al. (2019) Dataset:

- scATAC-seq data from human hematopoietic cells
- Published in Nature Biotechnology: "Droplet-based combinatorial indexing for massive-scale single-cell chromatin accessibility"

- Available through GEO: GSE123581
  (https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE123581)

6. ENCODE Project:

- The ENCODE project has various ATAC-seq datasets, including some single-cell data
- Available at: https://www.encodeproject.org/

7. Human Cell Atlas:

- Provides various single-cell datasets, including some scATAC-seq data
- Available at: https://data.humancellatlas.org/

8. Gonzalez-Blas et al. (2019) Dataset:
- scATAC-seq data from mouse brain development
- Published in Nature Communications: "Cis-regulatory dynamics during human development revealed by single-cell chromatin accessibility profiling"
- Available through GEO: GSE126074

These datasets cover a range of biological systems and experimental conditions, making them valuable resources for developing and testing scATAC-seq analysis methods. When using these datasets, always make sure to cite the original publications and adhere to any usage guidelines provided by the data generators.

# Experiment Environment

In most of the cases, the following environment is ensured for experiment purpose.

To set up an environment for experimenting with scATAC-seq data, we need a combination of software tools, programming languages, and computational resources. Here's a recommended setup:

1. Operating System:
   - Linux (Ubuntu, CentOS) or macOS
   - Windows with Windows Subsystem for Linux (WSL) can also work

2. Programming Languages:
   - R (version 4.0 or later)
   - Python (version 3.7 or later)

3. Integrated Development Environment (IDE):
   - RStudio for R
   - PyCharm or VS Code for Python
   - Jupyter Notebook/Lab for interactive analysis

4. Package Managers:
   - Conda (Miniconda or Anaconda)
   - BiocManager for Bioconductor packages in R

5. Essential R packages:
   - Seurat
   - Signac
   - ggplot2
   - GenomicRanges
   - chromVAR
   - cicero
   - ArchR

6. Essential Python packages:
   - scanpy
   - anndata
   - scikit-learn
   - numpy
   - pandas
   - matplotlib
   - seaborn

7. Bioinformatics tools:
   - Samtools
   - Bedtools
   - MACS2 (for peak calling)
   - Bowtie2 or BWA (for alignment, if working with raw data)

8. Version Control:
   - Git

9. Computational Resources:
   - A computer with at least 16GB RAM (32GB or more recommended for larger datasets)
   - Multi-core processor
   - Sufficient storage space (SSDs preferred for faster data access)

10. Cloud Platforms (optional, for larger datasets):
    - Amazon Web Services (AWS)
    - Google Cloud Platform

- Microsoft Azure

11. Containerization (optional):
   - Docker for creating reproducible environments

12. Workflow Management (optional):
   - Snakemake or Nextflow for building analysis pipelines

# Setup steps:

1. Install Miniconda and create a new environment:
   ```
   conda create -n scatac python=3.8 r-base=4.0
   conda activate scatac
   ```

2. Install R packages:
   ```R
   if (!requireNamespace("BiocManager", quietly = TRUE))
       install.packages("BiocManager")
   BiocManager::install(c("Signac", "Seurat", "GenomicRanges", "chromVAR", "cicero"))
   ```

3. Install Python packages:
   ```
   pip install scanpy anndata scikit-learn numpy pandas matplotlib seaborn
   ```

4. Install bioinformatics tools:
   ```
   conda install -c bioconda samtools bedtools macs2 bowtie2
   ```

5. Set up IDE and clone relevant repositories or download datasets.

Remember to regularly update packages and tools. This environment should provide a solid foundation for experimenting with scATAC-seq data, from preprocessing to advanced analyses.

# Good Learning Resources

Here are some excellent learning resources for getting started with scATAC-seq analysis:

1. Tutorials and Workshops:

  a) Signac Tutorial (R-based):
    - Official tutorial for the Signac package
    - https://docs.signac.io/en/latest/tutorial.html

  b) ArchR Tutorial (R-based):
    - Comprehensive tutorial for the ArchR package
    - https://www.archrproject.com/bookdown/

  c) Scanpy Tutorial (Python-based):
    - While primarily for scRNA-seq, it's useful for general single-cell analysis
    - https://scanpy-tutorials.readthedocs.io/en/latest/

  d) ENCODE ATAC-seq Workshop:
    - Covers ATAC-seq analysis, including single-cell
    - https://www.encodeproject.org/atac-seq/

2. Books and Comprehensive Guides:

  a) "Orchestrating Single-Cell Analysis with Bioconductor"
    - Includes chapters on ATAC-seq analysis
    - https://osca.bioconductor.org/

  b) "Current Best Practices in Single‑Cell ATAC‑Seq Analysis: A Tutorial"
    - Published in Molecular Systems Biology
    - https://www.embopress.org/doi/full/10.15252/msb.20188746

3. YouTube Channels and Videos:

  a) StatQuest with Josh Starmer:
    - Excellent explanations of statistical concepts in genomics
    - https://www.youtube.com/user/joshstarmer

  b) Bioinformagician Channel:
    - Covers various computational biology topics
    - https://www.youtube.com/@Bioinformagician

4. GitHub Repositories:

   a) Awesome Single Cell:
      - Curated list of single-cell resources, including ATAC-seq
      - https://github.com/seandavi/awesome-single-cell

   b) scATAC-pro:
      - Comprehensive scATAC-seq processing pipeline
      - https://github.com/wbaopaul/scATAC-pro

   c) learning-bioinformatics-at-home
      - resources for learning bioinformatics
      - https://github.com/harvardinformatics/learning-bioinformatics-at-home

5. Scientific Papers:

   a) "Single-cell ATAC sequencing analysis: From data preprocessing to hypothesis generation"
      - Comprehensive review of scATAC-seq analysis
      - https://pubmed.ncbi.nlm.nih.gov/32637041/

   b) "Computational Principles and Challenges in Single-Cell Data Integration"
      - Covers integration of scATAC-seq with other data types
      - https://www.nature.com/articles/s41587-021-00895-7

6. Community Forums:

   a) Bioconductor Support Forum:
      - https://support.bioconductor.org/

   b) Seurat GitHub Issues:
      - Often contains discussions on ATAC-seq analysis
      - https://github.com/satijalab/seurat/issues

# Sample Code on scATAC Data Analysis

We are going to show how anyone can work on scATAC-related data analysis with Python language. The whole experiment given below has been performed on an Apple Macbook M1 Air laptop.

## Importing necessary libraries

```
import scanpy as sc
import anndata as ad
import matplotlib.pyplot as plt
import numpy as np
import pandas as pd
import os
import scrublet as scr
```

This will import the necessary libraries for our research. If we get any error like a library is not found, then we can simply install the library using pip (pip install libraryname).

## Define dataset URL and filename (using a PBMC dataset as an example)

```
dataset_url =
"https://cf.10xgenomics.com/samples/cell-atac/1.1.0/atac_pbmc_10k_v1/atac_pbmc_10k_v1_filtered_peak_bc_matrix.h5"
dataset_filename = "atac_pbmc_10k_filtered_peak_bc_matrix.h5"
if not os.path.exists(dataset_filename):
    print(f"Downloading dataset from {dataset_url}...")
    os.system(f"wget -O {dataset_filename} {dataset_url}")
```

Here, we have added a small dataset from 10xgenomics and we are using the "atac_pbmc_10k_filtered_peak_bc_matrix.h5" file.

# Load the dataset from the provided URL or use a demo dataset from Scanpy

```
adata = sc.read_10x_h5(dataset_filename)
adata = sc.datasets.pbmc3k()
print(f"Loaded dataset with {adata.n_obs} cells and {adata.n_vars} genes.")
```

We are using the exact demo dataset that we have downloaded just now. It will load the dataset and print exactly how many cells and genes are in the dataset.

# Basic Quality Control and Filtering

```
mito_genes = adata.var_names.str.startswith('MT-')
adata.obs['percent_mt'] = np.sum(adata[:, mito_genes].X, axis=1) / np.sum(adata.X, axis=1)
adata.obs['total_counts'] = adata.X.sum(axis=1)
adata.obs['n_genes'] = (adata.X > 0).sum(axis=1)
adata.obs['log_counts'] = np.log1p(adata.obs['total_counts'])
sc.pl.violin(adata, ['n_genes', 'total_counts', 'percent_mt'], jitter=0.4, multi_panel=True)
sc.pl.scatter(adata, x='total_counts', y='n_genes')
sc.pl.scatter(adata, x='total_counts', y='percent_mt')
adata = adata[(adata.obs['n_genes'] > 200) & (adata.obs['n_genes'] < 2500) &
(adata.obs['percent_mt'] < 5), :]
print(f"After quality control filtering: {adata.n_obs} cells remain.")
```

We have performed quality control on single-cell RNA sequencing data by calculating and visualizing key metrics such as mitochondrial gene percentage and gene expression counts. We then filtered out low-quality cells based on these metrics using predefined thresholds. This process ensures that only high-quality cells are retained for further analysis, removing potential artifacts and improving the reliability of downstream results.

# Doublet Detection using Scrublet

```
scrub = scr.Scrublet(adata.X)  # Initialize Scrublet object
doublet_scores, predicted_doublets = scrub.scrub_doublets()
```

```
adata.obs['doublet_score'] = doublet_scores  # Add doublet scores to AnnData
adata.obs['predicted_doublet'] = predicted_doublets
scrub.plot_histogram()
plt.show()
adata = adata[~adata.obs['predicted_doublet'], :]
print(f"After doublet removal: {adata.n_obs} cells remain.")
```

We have used Scrublet to detect and remove potential doublets from our single-cell RNA sequencing data, calculating doublet scores and predictions for each cell. After visualizing the distribution of doublet scores and filtering out predicted doublets, we retained only high-quality singlets for further analysis.

# Normalization and Log Transformation

```
sc.pp.normalize_total(adata, target_sum=1e4)
sc.pp.log1p(adata)
```

We have normalized the data by scaling cell counts to a common total and applied a log transformation. This process standardizes the data across cells and prepares it for further analysis.

# Highly Variable Genes (HVGs) Identification

```
sc.pp.highly_variable_genes(adata, min_mean=0.0125, max_mean=3, min_disp=0.5)
sc.pl.highly_variable_genes(adata)  # Plot HVGs
adata = adata[:, adata.var.highly_variable]  # Filter to keep only HVGs
sc.pp.scale(adata, max_value=10)  # Scale each gene to unit variance and center around 0
```

We have identified and selected highly variable genes to focus on the most informative features in the dataset. We then scaled the selected genes to unit variance and centered them around zero, capping extreme values to mitigate outlier effects.

# Dimensionality Reduction using PCA and NMF

```
sc.tl.pca(adata, svd_solver='arpack')
sc.pl.pca_variance_ratio(adata, log=True)
```

We have performed Principal Component Analysis (PCA) on the data to reduce its dimensionality. We then plotted the explained variance ratio of the PCA components to visualize how much variance is captured by each component, helping to determine the optimal number of components for downstream analysis.

# Using Non-negative Matrix Factorization (NMF) for dimensionality reduction from sklearn.decomposition import NMF (optional)

```
if (adata.X < 0).any():
    print("Negative values detected. Adjusting to make all values non-negative...")
    adata.X = adata.X - np.min(adata.X) + 1
nmf_model = NMF(n_components=50, random_state=0)
adata.obsm['X_nmf'] = nmf_model.fit_transform(adata.X)
sc.pp.neighbors(adata, n_neighbors=15, use_rep='X_nmf')  sc.tl.umap(adata)
sc.pl.umap(adata, color=['total_counts', 'n_genes', 'percent_mt'])
```

We adjusted the data to ensure non-negative values and then applied Non-negative Matrix Factorization (NMF) for dimensionality reduction. Using the NMF-reduced data, we computed a neighborhood graph and generated a UMAP visualization, coloring cells by quality control metrics to examine potential patterns or biases in the reduced dimensional space.

## Clustering using the Leiden Algorithm

```
sc.tl.leiden(adata, resolution=0.8)
sc.pl.umap(adata, color='leiden')
sc.tl.tsne(adata, n_pcs=40)
sc.pl.tsne(adata, color='leiden')
```

We performed clustering on the data using the Leiden algorithm, which grouped similar cells together. We then visualized these clusters on both UMAP and t-SNE plots, providing two different perspectives on how the cells are organized in the reduced dimensional space.

## Advanced Differential Expression Analysis

```
sc.tl.rank_genes_groups(adata, 'leiden', method='wilcoxon')
sc.pl.rank_genes_groups(adata, n_genes=25, sharey=False)
sc.pl.rank_genes_groups_dotplot(adata, n_genes=4)
sc.pl.rank_genes_groups_matrixplot(adata, n_genes=4)
```

We conducted differential expression analysis using the Wilcoxon rank-sum test to identify marker genes for each cluster. We then visualized the results using various plots, including a ranking of top marker genes, a dotplot, and a matrixplot, to better understand the gene expression patterns that characterize each cluster.

## Save Processed and Improved Data

```
output_file_improved = "pbmc3k_improved_data_v2.h5ad"
adata.write(output_file_improved)
print(f"Enhanced dataset saved to {output_file_improved}.")
```

We saved the processed and analyzed single-cell data, including all results and annotations, to a file in the AnnData (.h5ad) format. This step ensures that our enhanced dataset is preserved for future use or sharing with other researchers.

# Conclusion

In conclusion, we have performed a comprehensive analysis of single-cell RNA sequencing data, including quality control, doublet detection, normalization, feature selection, dimensionality reduction, clustering, and differential expression analysis. This pipeline has transformed raw cellular data into biologically meaningful clusters and identified key marker genes, providing valuable insights into cell populations and their characteristics. The final saved dataset encapsulates all these analytical steps, serving as a robust foundation for further biological interpretation and downstream analyses.