# Clustering Human Organ Cells

## CS-297

Name: Swathi M.V.S
ID: 016669757

# Outline

# Dataset

H5ad Data Format

- Specific to R language for single-cell RNA sequencing
- Widely adapted and used persistent on-disk storage
- H5ad format - based on standard h5 format - Hierarchical Data Formats (HDF)
- HDFs - use multidimensional arrays to store large amounts of data
- H5 format - used to store scientific data that is well-organized for quick retrieval and analysis
- Developed by Theislab with extensive support in Python

# Dataset

Tabula Sapiens - Heart

- Specific to R language.
- The H5ad format is based on the standard h5 format, a Hierarchical Data Formats (HDF) used to store large amounts of data in the form of multidimensional arrays. The H5 format is primarily used to store scientific data that is well-organized for quick retrieval and analysis.

# Dataset

Tabula Sapiens - Heart

- The heart dataset contains 11,505 cells and 58,604 genes

```
[ ]  # Get the dimensions of the data
     print("Number of Cells:", adata.n_obs)
     print("Number of Genes:", adata.n_vars)

     Number of Cells: 11505
     Number of Genes: 58604
```

# Libraries

- Scanpy - It's commonly used in bioinformatics for analyzing single-cell genomics data.
- Anndata – It is used to to store and manipulate data associated with individual cells and handling annotated data
- Matplotlib - Library for creating static, animated, and interactive visualizations
- Seaborn - Built on top of Matplotlib and provides a high-level interface for creating informative and attractive statistical graphics
- Pandas - Data manipulation and analysis library. It provides data structures like DataFrame and Series
- Numpy - fundamental package for scientific computing in Python
- Os – module to interact with the operating system

# Attributes

- The data contains the following attributes

```
AnnData object with n_obs × n_vars = 11505 × 58604 obs:
'assay_ontology_term_id', 'donor_id', 'anatomical_information',
'n_counts_UMIs', 'n_genes', 'cell_ontology_class',
'free_annotation', 'manually_annotated', 'compartment',
'sex_ontology_term_id', 'disease_ontology_term_id',
'is_primary_data', 'organism_ontology_term_id', 'suspension_type',
'cell_type_ontology_term_id', 'tissue_ontology_term_id',
'development_stage_ontology_term_id',
'self_reported_ethnicity_ontology_term_id', 'cell_type', 'assay',
'disease', 'organism', 'sex', 'tissue', 'self_reported_ethnicity',
'development_stage' var: 'feature_type', 'highly_variable',
'means', 'dispersions', 'dispersions_norm', 'mean', 'std',
'ensembl_version', 'feature_is_filtered', 'feature_name',
'feature_reference', 'feature_biotype' uns: '_scvi',
'_training_mode', 'assay_colors', 'cell_ontology_class_colors',
'dendrogram_cell_type_tissue',
'dendrogram_computational_compartment_assignment',
'dendrogram_consensus_prediction', 'dendrogram_tissue_cell_type',
'donor_id_colors', 'hvg', 'neighbors', 'schema_version',
'sex_colors', 'tissue_colors', 'title', 'umap' obsm: 'X_pca',
'X_scvi', 'X_scvi_umap', 'X_umap' obsp: 'connectivities',
'distances'
```

# Attributes

- The data contains the following attributes

```
print(adata)
```

```
AnnData object with n_obs x n_vars = 11505 x 58604
    obs: 'assay_ontology_term_id', 'donor_id', 'anatomical_information', 'n_counts_UMIs', 'n_genes', 'cell_ontology_class', 'free_annotation', 'manually_annotated', 'compartment', 'sex
    var: 'feature_type', 'highly_variable', 'means', 'dispersions', 'dispersions_norm', 'mean', 'std', 'ensembl_version', 'feature_is_filtered', 'feature_name', 'feature_reference', 'f
    uns: '_scvi', '_training_mode', 'assay_colors', 'cell_ontology_class_colors', 'dendrogram_cell_type_tissue', 'dendrogram_computational_compartment_assignment', 'dendrogram_consensu
    obsm: 'X_pca', 'X_scvi', 'X_scvi_umap', 'X_umap'
    obsp: 'connectivities', 'distances'
```

# Attributes

- Gene expression - the process by which the information encoded in a gene is turned into a function
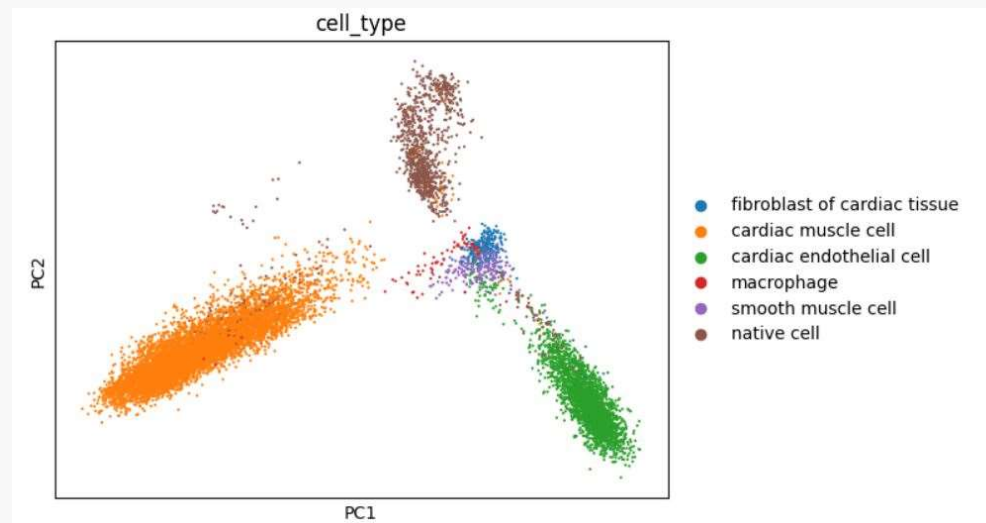
# Dimensionality Reduction

- UMAP (Uniform Manifold Approximation and Projection) - fairly flexible non-linear dimension reduction algorithm

# Dimensionality Reduction

- PCA - fairly flexible non-linear dimension reduction algorithm

# Dimensionality Reduction