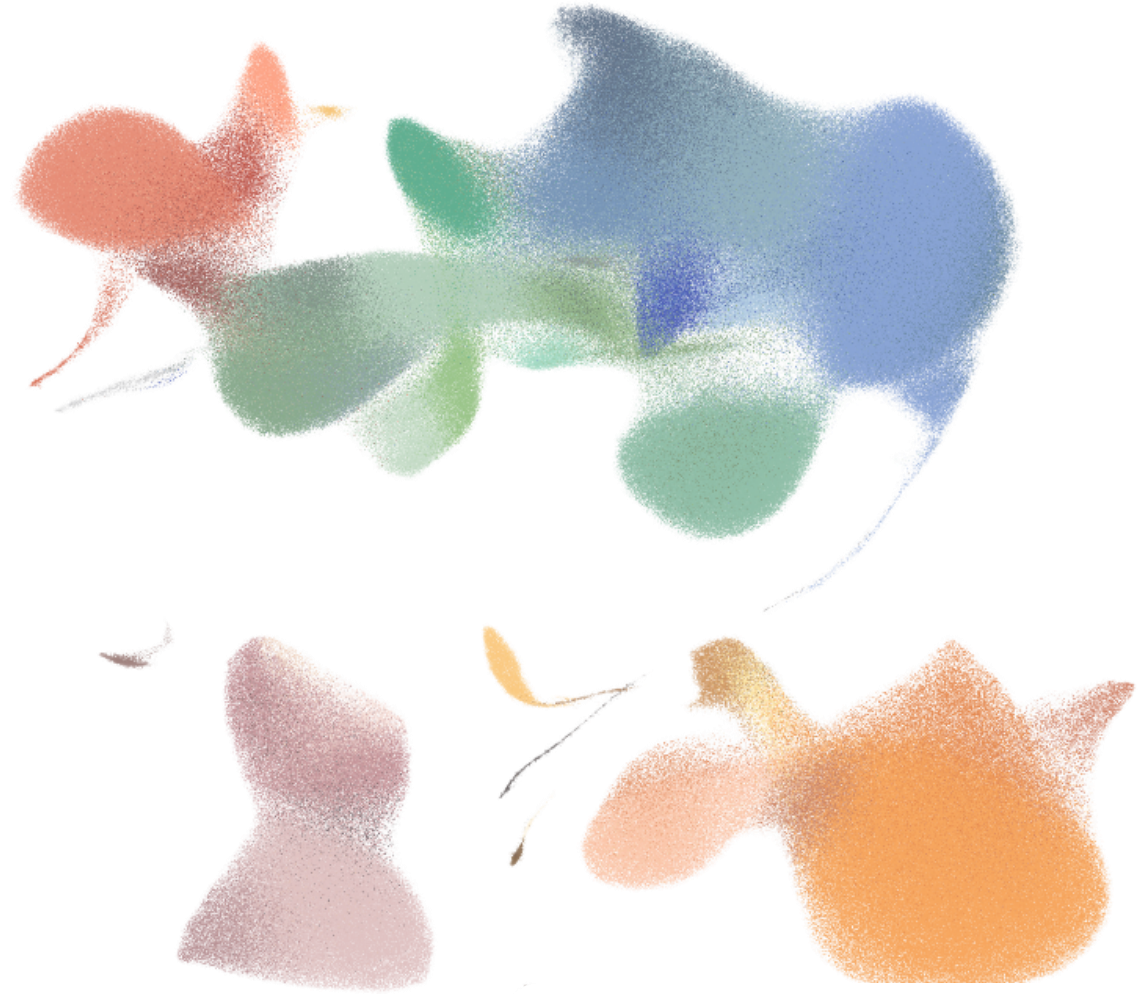# Clustering
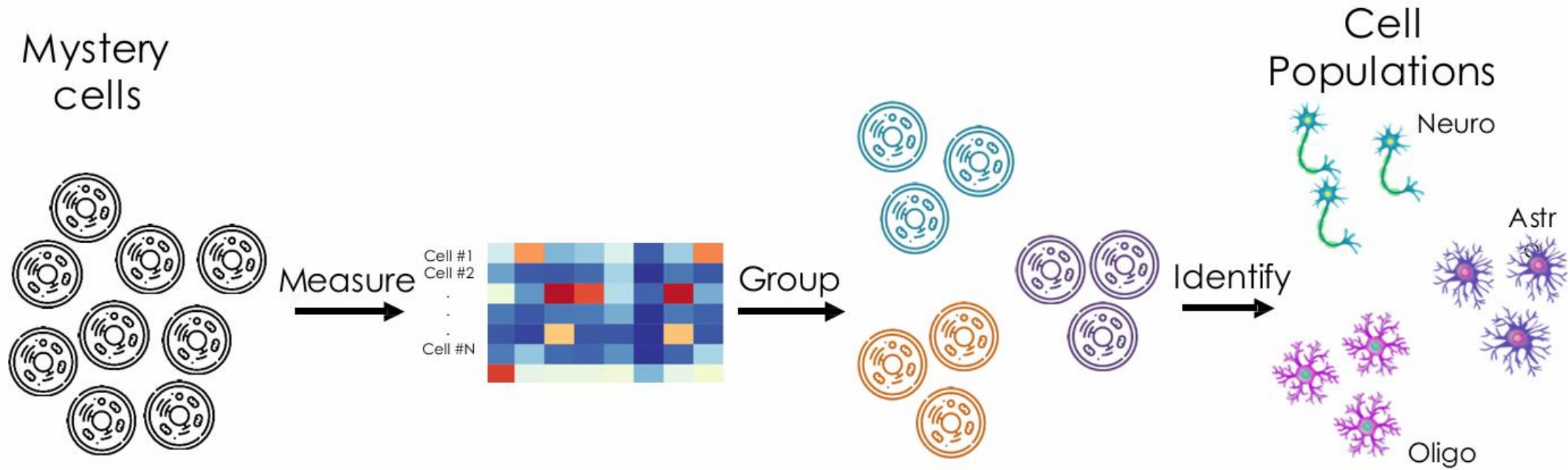
*Yusuf Caglar Odabasi*

*December 1.-3. 2025*

**Course on scRNA-seq Data Analysis**
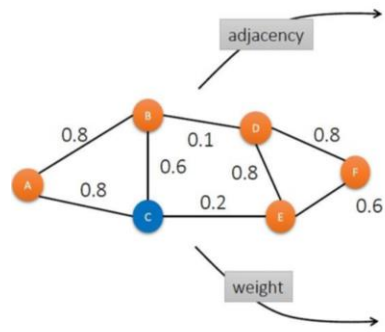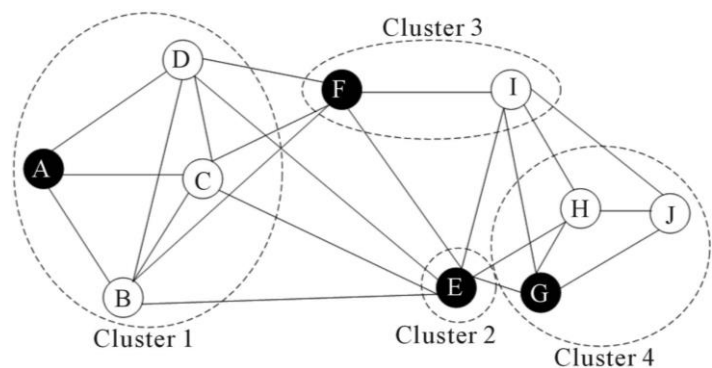
The goal of clustering is to group cells with similar gene expression profiles

# Concept of clustering

**Graph-based clustering**



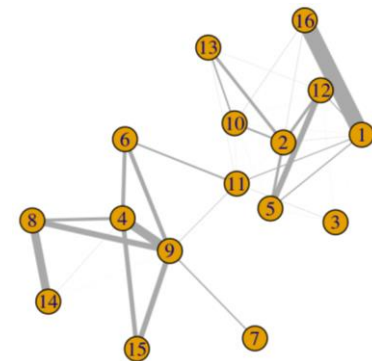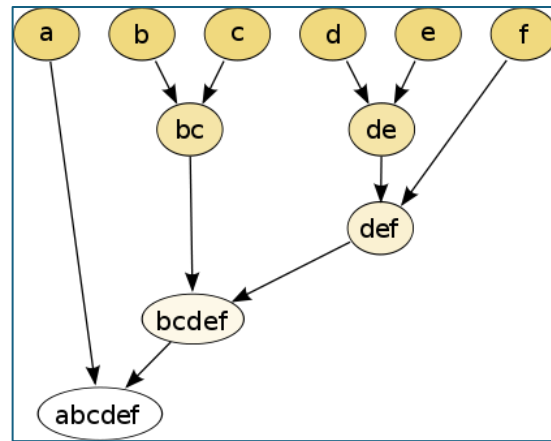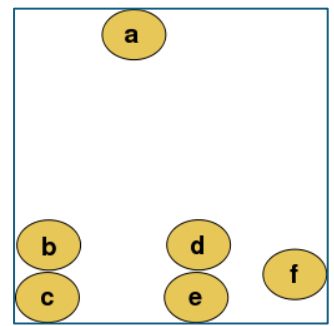**Hierarchical clustering**



*K*-means clustering

# Graph-based clustering

**-** K-nearest neighbour (KNN) graph based on the euclidean distance in PCA space.

Two vertices p and q are connected by an edge, if the distance between p and q is among the k-th smallest distances from p to other nodes.
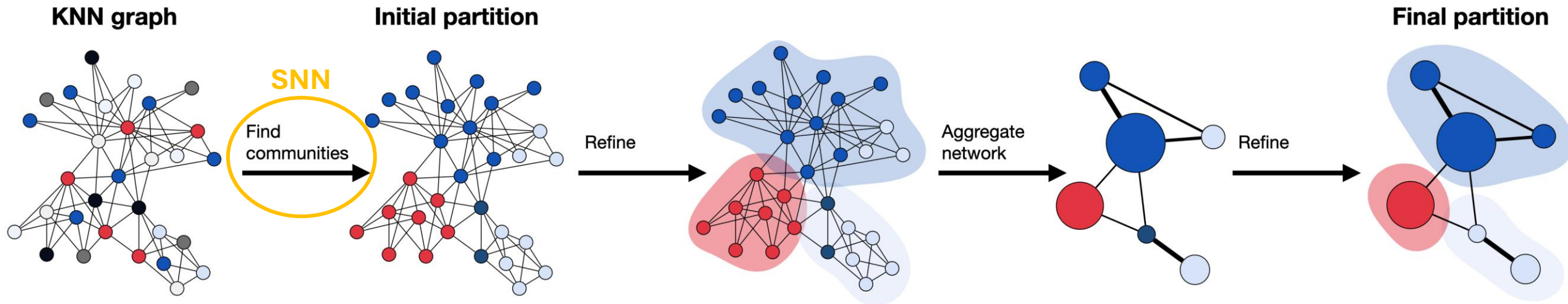
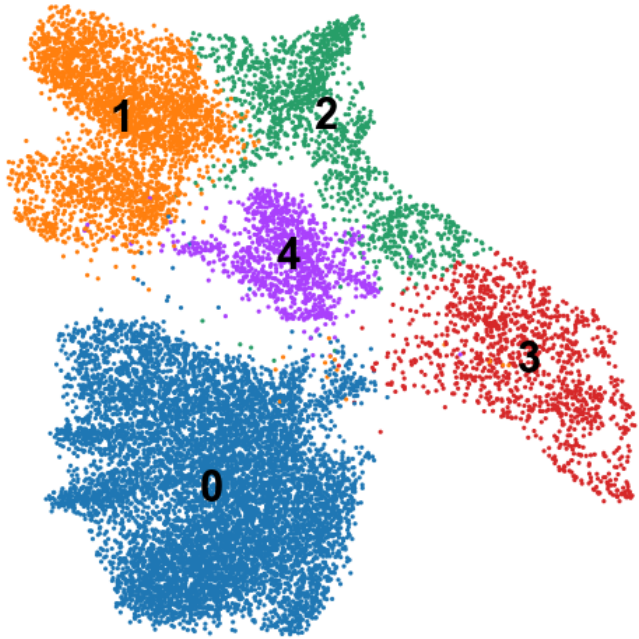**-** Shared-nearest neighbour (SNN) graph

For each pair of cells (nodes), the number of shared neighbours is counted (according to the KNN graph). An edge is created between two cells if they share a sufficient number of nearest neighbours (above a certain threshold).
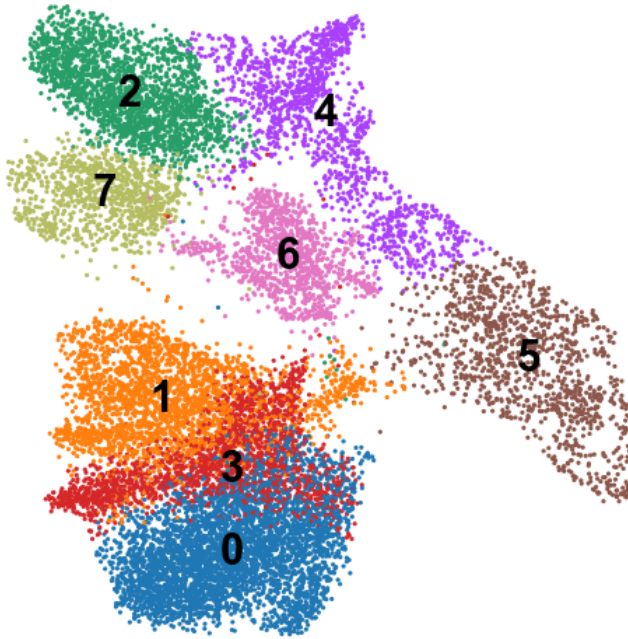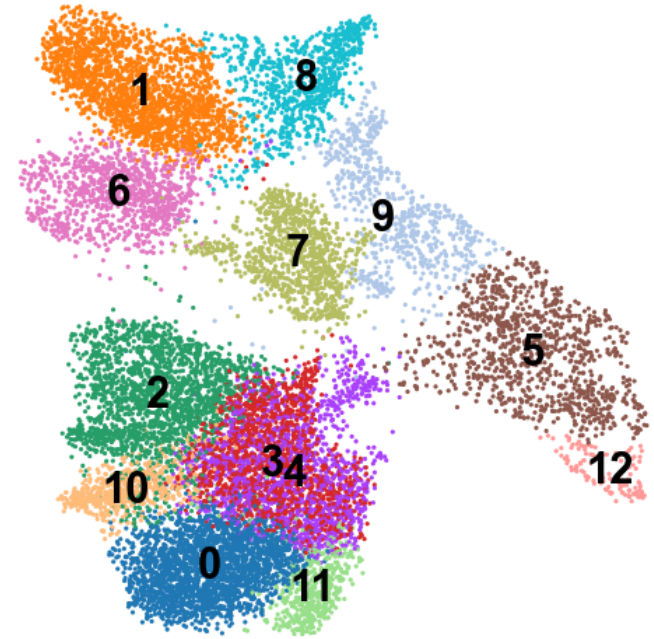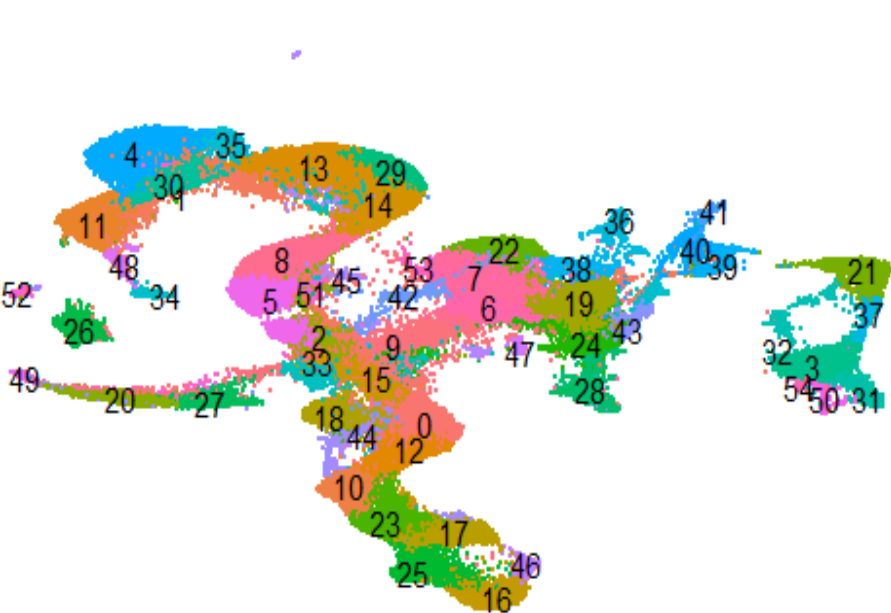


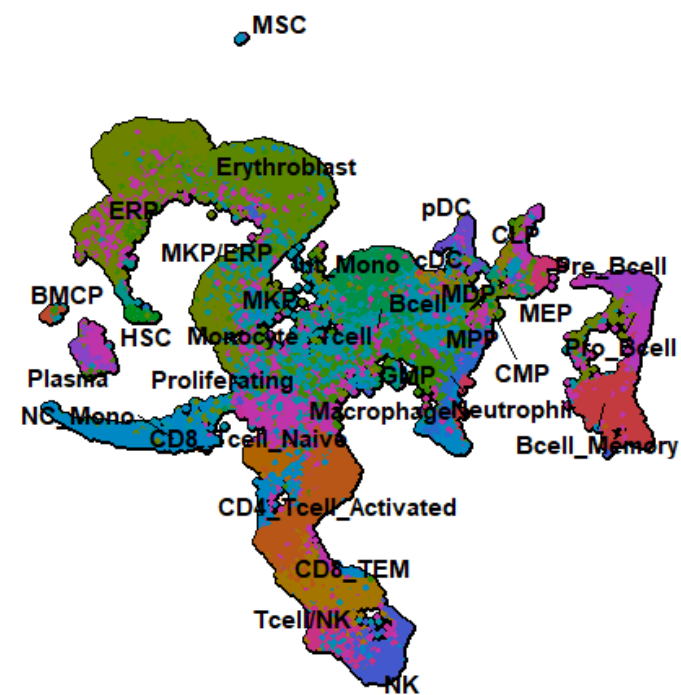**Leiden and Louvain**

leiden_res0_25   leiden_res0_5   leiden_res1

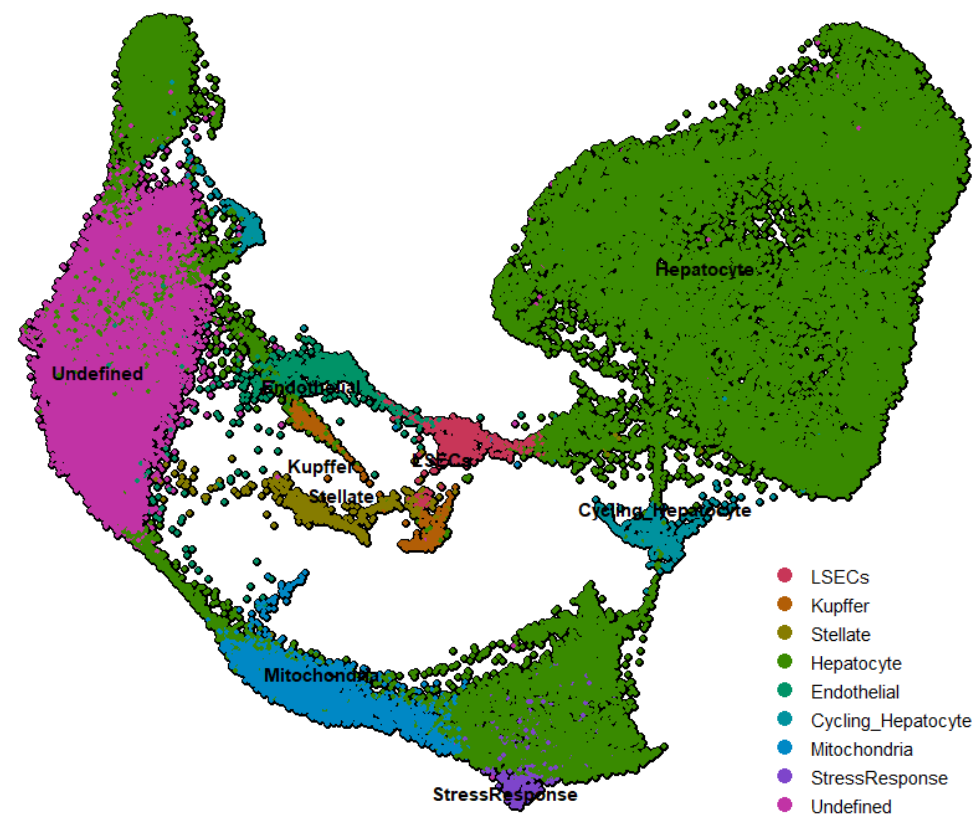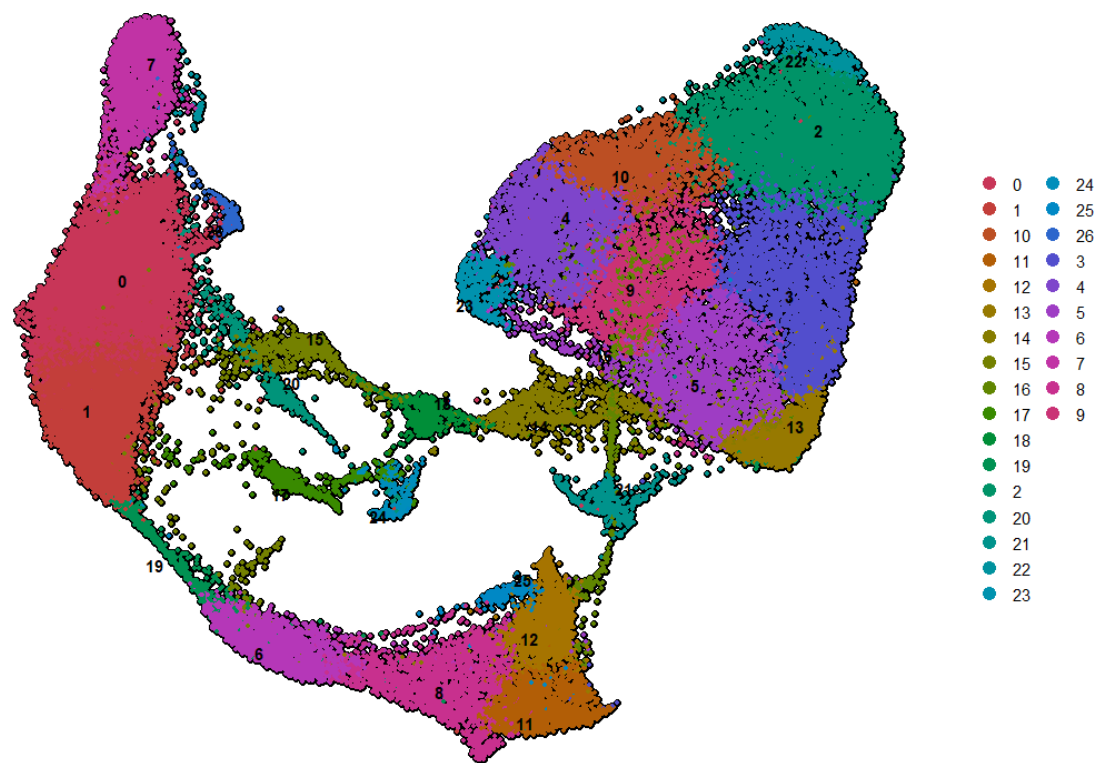There is not a correct number of clusters, it will depend on the context and biological question

Legend (sample identities):
- BOD17-1
- BOD17-3
- BOD18KV
- BOD19-1
- BOD20-3
- BOD21-3
- BOD22-1
- OBE-DIA 1-1
- OBE-DIA 2-1

Cluster legend:
0, 1, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 2, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 3, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 4, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 5, 50, 51, 52, 53, 54, 6, 7, 8, 9

Cell type legend:
- Bcell
- Bcell_Memory
- BMCP
- CD4_Tcell_Activated
- CD8_Tcell_Naive
- CD8_TEM
- cDC
- CLP
- CMP
- ERP
- Erythroblast
- GMP
- HSC
- Int_Mono
- Macrophage
- MDP
- MEP
- MKP
- MKP/ERP
- Monocyte
- MPP
- MSC
- NC_Mono
- Neutrophil
- NK
- pDC
- Plasma
- Pre_Bcell
- Pro_Bcell
- Proliferating
- Tcell
- Tcell/NK
- Transitional_Bcell

- https://satijalab.org/seurat/articles/integration_introduction

- https://www.sc-best-practices.org/cellular_structure/integration.html

- https://www.singlecellcourse.org/biological-analysis.html#clustering-introduction

- https://bioconductor.org/books/3.12/OSCA/clustering.html#k-means-clustering

- https://github.com/quadbio/scRNAseq_analysis_vignette/blob/master/Tutorial.md#step-2-3-data-integration-using-liger