# KSBi-BIML 2026

**Bioinformatics & Machine Learning(BIML) Workshop for Life Scientists**

생명정보학 & 머신러닝 워크샵 (온라인)

# Single-cell RNA-sequencing analysis: Assignment of cell types

김규태 _ 아주대학교

# KSBi-BIML 2026

## Bioinformatics & Machine Learning (BIML) Workshop for Life Scientists

한국생명정보학회가 주최하는 BIML-2026 동계 Bioinformatics &Machine Learning 교육 워크숍에 여러분을 초대합니다.

BIML 워크숍은 생명정보학 연구자들이 최신 AI바이오 분야의 인공지능 기반 분석 기술과 바이오 데이터 분석 기법을 이론과 실습을 통해 체계적으로 배울 수 있는 전문 교육 프로그램입니다. 2015년에 시작된 BIML 워크숍은 올해로 12년 차를 맞이하며, 국내 생명정보학 분야의 최초이자 최고 수준의 교육 프로그램으로 자리 잡았습니다. 이번 워크숍은 크게 인공지능바이오(AI바이오) 분야와 디지털바이오 분야, 두 분야로 구성됩니다.

AI바이오 분야에서는 생명정보 분석에 폭넓게 응용되고 있는 다양한 인공지능 기반 자료 모델링 기법을 다룰 예정입니다. 특히, 인공지능 심층학습을 활용한 단백질 구조 예측, 유전체 분석, 신약 개발에 대한 이론 및 실습 강의가 진행됩니다.

또한 디지털바이오 분야에서는 단일세포오믹스, 공간오믹스, 멀티오믹스, 메타오믹스에 대한 강의도 마련되어 있어, 연구자들의 분석 역량 강화에 실질적인 도움을 줄 것으로 기대됩니다.

또한 2024년부터 추가된 의료정보 자료 분석을 다루는 강의를 올해도 지속해서 운영하고자 합니다. 이는 최근 의료정보 자료 분석에 관한 연구 수요 증가를 반영한 것으로, 관련 연구를 수행하는 의과학자 및 의료정보 연구자들에게 유용한 지침을 제공할 것입니다.

또한, 올해도 생명정보학 기술의 다양화에 발맞춰 온라인 강좌를 대폭 확대했습니다. 올해는 무료 강좌 10개를 포함한 총 40개 이상의 강좌가 개설되며, 연구 주제에 맞는 강좌 추천과 강연료 할인 혜택도 제공합니다.

BIML-2026는 국내 주요 연구 중심 대학의 전임 교수 및 각 분야 최고 전문가들의 강의로 구성되어 있으며, 기초 이론부터 최신 연구 동향까지 아우르는 심도 있는 교육의 장이 될 것으로 확신합니다.

여러분의 많은 관심과 참여를 기대합니다!

2026년 2월

**한국생명정보학회장 류 성 호**

# Single-cell RNA-sequencing analysis: Assignment of cell types

본 강의는 단일세포 전사체 데이터 분석의 기본적인 측면을 다룬다. 단일세포 수준으로 분석하는 것이 왜 중요한지에 대한 개론을 제공하며, 데이터 유형의 구조와 형식을 설명하고, 데이터 전처리 과정을 이해할 수 있도록 이론과 함께 실습 강의를 제공한다. 또한, 단일세포 전사체 데이터를 이용한 세포 유형을 결정하는 전반적인 과정을 이해할 수 있다. 이를 통해 학습자들은 단일세포 연구에서 데이터를 처리하고 세포 유형을 파악하는데 필요한 기초적인 지식을 습득하게 된다.

강의는 다음의 내용을 포함한다:

- 단일세포 전사체 데이터 분석의 중요성과 의의를 이해
- 단일세포 전사체 데이터의 구조와 형식에 대해 학습
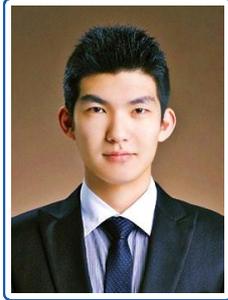- 단일세포 전사체 데이터를 활용하여 세포 유형을 할당하는 과정을 이해

\* 교육생준비물:

Rstudio 및 Seurat (R package)가 설치된 노트북 (메모리 8GB 이상, 디스크 여유공간 30GB 이상)

\* 강의 난이도: 초급

\* 강의: 김규태 교수 (아주대학교의과대학 생리학교실)

# Curriculum Vitae

## Speaker Name: Kyu-Tae Kim, Ph.D.

▶ **Personal Info**

| | |
|---|---|
| Name | Kyu-Tae Kim |
| Title | Assistant Professor |
| Affiliation | Ajou University School of Medicine |

▶ **Contact Information**

| | |
|---|---|
| Address | 164, Wolrd cup-ro, Yeongtong-gu, Suwon 16499 |
| Email | kimqtae@ajou.ac.kr |

---

**Research Interest**

Immunogenomics, Cancer evolution, Computational Biology

**Educational Experience**

| | |
|---|---|
| 2010 | B.S., Konkuk University, Seoul, Korea |
| 2012 | M.S., Seoul National University, Seoul, Korea |
| 2015 | Ph.D., Seoul National University, Seoul, Korea |

**Professional Experience**

| | |
|---|---|
| 2013-2017 | Researcher, Samsung Genome Institute, Samsung Medical Center, Seoul, Korea |
| 2017-2019 | Postdoctoral Fellow, New York Genome Center, NY, USA |
| 2020- | Assistant Professor, Ajou University School of Medicine, Suwon, Korea |

**Selected Publications (5 maximum)**

1. Determinants of Response and Intrinsic Resistance to PD-1 Blockade in Microsatellite Instability-High Gastric Cancer, Cancer Discovery, 2021 (corresponding author)

2. Somatic mutations and cell identity linked by Genotyping of Transcriptomes, Nature, 2019 (first author)

3. SIDR: simultaneous isolation and parallel sequencing of genomic DNA and total RNA from single cells, Genome Research, 2018 (first author)

4. Application of single-cell RNA sequencing in optimizing a combinatorial therapeutic strategy in metastatic renal cell carcinoma, Genome Biology, 2016 (first author)

5. Single-cell mRNA sequencing identifies subclonal heterogeneity in anti-cancer drug responses of lung adenocarcinoma cells, Genome Biology, 2015 (first author)

# KSBi-BIML 2024

## Single-cell RNA-sequencing analysis: Assignment of cell types (part1)

Kyu-Tae Kim
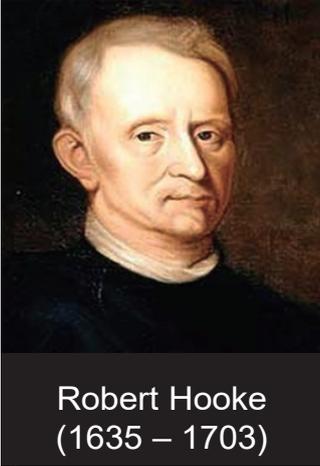Ajou University School of Medicine

---
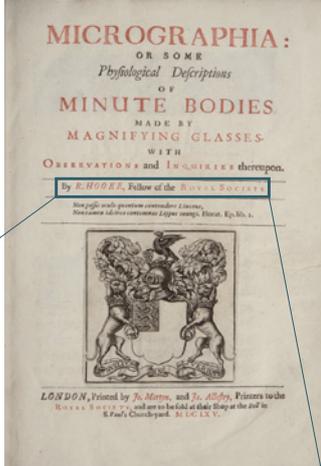
## 본 교육의 목표와 특징

### 단일세포 전사체 데이터 전분석

- 단일세포 전사체 데이터 분석의 의의를 이해한다.

- 단일세포 전사체 데이터의 구조와 형식을 이해한다.

- 단일세포 전사체 데이터의 전분석 과정을 이해한다.

- 단일세포 전사체 데이터 normalization 과정을 이해한다.

- 단일세포 전사체 데이터 batch 제거 과정을 이해한다.

# Cell: The basic unit of life

Drawing by Hooke



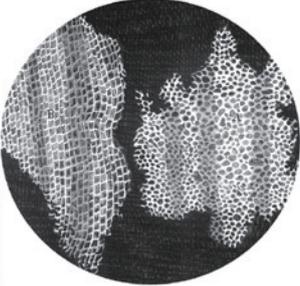Robert Hooke was the first
to apply the word 'Cell' to biological objects (Cork).



Robert Hooke
(1635 – 1703)

MICROGRAPHIA:
OR SOME
Physiological Descriptions
OF
MINUTE BODIES
MADE BY
MAGNIFYING GLASSES.
WITH
OBSERVATIONS and INQUIRIES thereupon.

By R. HOOKE, Fellow of the ROYAL SOCIETY.

By R. HOOKE, Fellow of the ROYAL SOCIETY.

Captured picture of Cork tissue



3



Cellular heterogeneity

# Cell: The basic unit of life

## From a cell, then ...

**Pluripotent cell**

**Dynamic acetylation/ phosphorylation**

**Increase in methylation**

**Cell differentiation**

Cell state A    Cell state B    Cell state C    Cell state D

**Waddington's model**

---

# Cell: The basic unit of life

**Tumor Micro-Environment**

Within-cell-type differences

Blood vessel

f

a b
c
g
d
e

Lymphatic vessel

Active tumor area

Tumor necrotic area

| | |
|---|---|
| Capillaries | |
| RBC | |
| Macrophage | |
| T lymphocyte | |
| B lymphocyte | |
| Normal fibroblasts | |
| endothelial cells | |
| pericytes | |
| Normal cells | |
| Malignant cells in necrotic or hypoxic area | |
| Cancer associated fibroblasts (CAFs) | |
| Dendritic cell | |
| Adipocyte | |
| Malignant cells | |
| NK and NKT cells | |

# Why single-cell sequencing?

Bulk analysis vs. Single-cell RNA-seq

---

**"No! Sometimes the Sum of the Parts (single-cells) is Greater than the Whole (bulk)."**
(original phrase by Aristotle, "The Whole is Greater than the Sum of its Parts.")



Variant allele frequency
or
Level of gene expression

# The bulk measurement is
## the stochastic average of single cells



Kim KT, Lee HW, Lee HO et al., 2015 *Genome Biol.*

# Single-cell analysis – a brief history

'Single-cell sequencing'
Methods of the Year 2013

Rapid progress in
'Single-cell sequencing'

[Tracking development cell by cell]
Breakthrough of the Year
2018 Science



adapted from Wang *et al. Mol Cell* 2015



10

# Single-cell analysis – a brief history



**Above timeline (top labels, left to right):**

- Single cell mRNA sequenced, *Tang, et.al.*
- First single-cell mRNA sequenced, *Tang, et.al.*
- First single-cell exome sequenced, *Xu, et.al.*
- SC RNA-Seq of immune cells, *Shalek, et.al.*
- SC ATAC-Seq of human cell lines, *Cusanovich, et.al.*
- Single cell qPCR of neurons, *Eberwine, et.al.*
- DNA-Seq of single human cancer cell, *Navin, et.al.*
- Single cell WGS of Neurons, *Evrony, et.al.*
- First single-cell T-cell epigenome with HiC, *Nagano, et.al.*
- Organ lineage tracing with SC RNA-Seq, *Truetlein, et.al.*
- Blastocyst E3.5 (ICM, TE)

**Timeline years:** 1992 — 2003 — 2009 — 2010 — 2011 — 2012 — 2013 — 2014 — 2015

**Below timeline (bottom labels, left to right):**

- SC whole transcriptome microarrays
- STRT-Seq 5' end Seq *Islam, et. al.*
- Fluidigm C1 AutoPrep released
- CEL-Seq Remove PCR bias with IVT *Hashimshony, et. al.*
- MALBAC *Zong, et.al.*
- Incorporation of UMI's to methods *Jaitin, et. al. Islam, et. al.*
- SMART-Seq Full length pre-amp *Ramskold, et. al.*
- SMART-Seq2 Increased sensitivity *Picelli, et.al.*
- Drop-Seq *Macosko, et. al.*

Kolodziejczyk, et. al. 2015; Navin, 2015

11

---

# Trends: Increasing Dimensionality & More Cells



Sarah Teichmann group, 2018, *Nat Proc*.

12

# [Experimental Approaches]



Multi-modal profiling methods at single-cell resolution

- DNA-seq + RNA-seq = **SIDR-seq**, **G&T-seq**, **DR-seq**
- DNA-seq + RNA-seq + Methyl-seq = **Trio-seq**
- RNA-seq + ATAC-seq = **sciCAR**
- RNA-seq + TCR/BCR = **(10X) 5'GEX with Immune Cell profiling**
- Epitope-profiling + RNA-seq + = **CITE-seq**
- Genotyping + RNA-seq = **GoT**
- Genetic screening with CRISPR + RNA-seq = **Perturb-seq**
- and……

13

# [Computational Approaches]



Farrell JA, Wang Y et al., 2018 *Science*

Manno GL et al., 2018 *Nature*

14

**[Experimental & Computational Approaches]**

genotyping of singular/multiple mutations per cell

GoT (linear/circ)

profiling of epigenetic mutations

Methylation (RRBS)

ATAC (10X)

profiling of cis-regulatory modes

Integrated analysis by matching cell barcodes

profiling of transcriptional readouts

RNA (10X)

long-read (ONT)

reliably detecting of alternative splicing with high-throughput

Clinical features

15



**Highly Dimensional Single-cell Data Sets**

**# Cells** x **# Features** x **# Time Points** x **# Technologies**

dissected by

Sophisticated Analytical Design with Massive Computational Power

16

# Basic single-cell analysis workflow

**Isolating Single Cells from Heterogeneous Population of Tumor**

**Amplifying Genome from Single Cells for Sequencing**

**Taking a Closer Look at Individual Cancer Cells**

- • Micropipetting
- • Laser capture microdissection
- • FACS
- • Microfluidic circuits
- • Droplet-based microfluidics

(DNA)
- • MALBAC
- • MDA
- • LIANTI

(RNA)
- • STRT-seq
- • CEL-seq
- • SMART/SMART2/SMART3-seq
- • Droplet-based amplification (Drop-seq, inDrop, 10X)

(statistical/algorithmical mining)

17

---

# At the initial stage of single-cell field

**Isolating Single Cells from Heterogeneous Population of Tumor**

**C₁™ Single-Cell AutoPrep System**
Fluidigm

**Amplifying Transcriptome from Single Cells for Sequencing**

**High-throughput sequencing**
illumina
RNA sequencing

**96** wells

Dead or live?
Singlet or doublet?

**40–80** live cells captured

How many genes detected?
How many reads aligned?
Mitochondria fraction?

**20–60** cells of QC-passed

18

# Basic data processing workflow

# Basic data processing workflow

for full-length/high-depth of several single-cells



Kim KT, Lee HW, Lee HO et al., 2015 *Genome Biol.*

# Trends: Increasing Dimensionality & More Cells



Sarah Teichmann group, 2018, *Nat Proc*.

# Basic single-cell analysis workflow

# Pre-processing pipeline: 10X  CellRanger



```
# /data/users/kimqt2/Projects/chonh_covid19/run_CellRanger.sh
/data/users/kimqt2/program/cellranger-3.1.0/cellranger count \
--id=20_00028_LI_SING \
--fastqs=/data/users/kimqt2/Projects/chonh_covid19/Lung_Fastq/ \
--transcriptome=/data/users/kimqt2/ref/tenX/refdata-cellranger-GRCh38-3.0.0_withSARS_COV2_SNU01 \
--expect-cells=5000 \
--localcores=30 \
--localmem=32
```

--> Output: Gene-level expression matrix per cell

---

# CellRanger (10X Genomics)

**1. Read Trimming**

> Detection/trimming of technically-induced sequence (TSO, template switch oligo)

**2. Read Alignment**

> Splicing-aware alignment of cDNA sequences to the genome reference using STAR

**3. Calling cell barcodes and UMI**

> Error-aware statistical correction of barcodes and UMI

**4. Basic subclustering and dimensional reduction**

# Identifying error-corrected barcode sequence



Whitelisted barcodes or Not?

whitelist barcodes
(n = ~737,000)

(for this comparison, no mismatches allowed)

737K-august-2016.txt
provided from 10x

- All reads are perfectly matched with the whitelist barcodes except for r0, r11, r15
- Reads of r0, r11, r15 are tested for possibility of replacing barcodes, in the next step

# Identifying error-corrected barcode sequence



whitelist barcodes

wb1
wb2  ACGTTAACGATACAAA
wb3  ACGTTAACGATAAAAA
wb4  ACGATAAGGATAAAAA

wb1005  ACGATAACGATAACAAA
wb1006

wb737280

Hamm=1 ?

ACGATAACGATAAAAA

r0

r11

r15

(1) Candidate barcodes of replacement ← Hamming−distance =1 comparing with the whitelist barcodes

- r11: only one base at any position is different with wb3, wb4, wb1005, but, Hamming=2 with wb2
  > r11 has three candidate barcodes

- r15: two candidate barcodes of wb4, wb1006

- r0: No Hamm=1

- Estimate probabilities to be changeable or not with r11, r15 that have candidate barcodes, in the next step

# Identifying error-corrected barcode sequence

| candidate barcodes | r11 | | | r15 | |
|---|---|---|---|---|---|
| | wb3 | wb4 | wb1005 | wb4 | wb1006 |
| Posterior probability | 0.9976318 | 0.002368154 | 0 | 1 | 0 |

Statistically Reliable to be replaced?



whitelist barcodes
wb1
wb2
wb3 ACGTTAACGATAAAA
wb4 ACGATAAGGATAAAAA
wb1005
wb1006
wb737280

Posterior probability

wb3 0.997631  wb4 0.002368  wb1005 0

A  C  G  A  T  A  A  C  G  A  T  A  A  A  A

r11
r15

(3) supporting reads + base quality

> posterior probability >=0.99 --> Barcode replacement!
– r11: among candidate barcodes of wb3, wb4, wb1005, wb3 has the highest posterior probability (= 0.9976318)
  > if this maximum posterior probability >=0.99, then replace the observed barcode with this wb3
– r15: between candidate barcodes of wb4, wb1006, wb4 has the maximum posterior probability (= 1)
  > if this maximum posterior probability >=0.99, then replace the observed barcode with this wb4
– If the maximum posterior probability <0.99, regarded as not significant to be replaced of barcode and rescued of the read

r11          replacing barcode          wb3 ... r11
r15                                      wb4 ... r15

(actually observed barcode in r11)          →  wb3
(actually observed barcode in r15)          →  wb4

# Identifying error-corrected barcode sequence



Unique total barcodes

in white-list #37696
out of white-list #149609

Replacing barcodes

in white-list #37696 replaced
out of white-list & not replaced #131006
out of white-list #18603



Number of reads

Priming 90.66%

out of white-list

replaced out of white-list

in white-list          in white-list

A. Identification of reads that have priming          B. Replacing barcodes

# Output BAM

(White-listed barcodes)
- 10X v2.chemistry: `737K-august-2016.txt`
- 10X v3.chemistry: `3M-february-2018.txt`

(Length of UMI)
- 10X v2.chemistry: `10`
- 10X v3.chemistry: `12`

# Output matrices



← sparse matrices for gene expressi

# Estimation of relative level of gene expression & Normalization of their abundances



duplication을 구별하기 위해서 Unique sequence tagging: UMI (Unique Molecular Identifier)

# Estimation of relative level of gene expression & Normalization of their abundances



reads with replaced CB

([r1, r2] , [r8, r9, r10], [r13, r14,r15] <--- identical 'barcode+UMI' (duplicate)

# Output matrices

---

## <span style="color:red">Seurat R package:</span>
## the most popular tool package processing 10X output data

https://satijalab.org/seurat

1) Read the 10X output data

2) QC and select cells for further analysis

3) Normalize the data

4) Detection of variable genes across the single cells

5) Scale the data and remove unwanted sources of variation

6) Perform linear dimensional reduction

7) Determine statistically significant principal components

8) Cluster the cells

9) Run non-linear dimensional reduction

10) Find differentially expressed genes (cluster biomarkers)

11) Assign cell type identity to clusters

12) Further sub-dissect within cell types

# Data loading (practice)

```
./cellrangers/
├── ctl.1
│   ├── barcodes.tsv.gz
│   ├── features.tsv.gz
│   └── matrix.mtx.gz
├── ctl.2
│   ├── barcodes.tsv.gz
│   ├── features.tsv.gz
│   └── matrix.mtx.gz
├── ctl.3
│   ├── barcodes.tsv.gz
│   ├── features.tsv.gz
│   └── matrix.mtx.gz
├── luad.1
│   ├── barcodes.tsv.gz
│   ├── features.tsv.gz
│   └── matrix.mtx.gz
├── luad.2
│   ├── barcodes.tsv.gz
│   ├── features.tsv.gz
│   └── matrix.mtx.gz
└── luad.3
    ├── barcodes.tsv.gz
    ├── features.tsv.gz
    └── matrix.mtx.gz

6 directories, 18 files
```

# Data loading

```r
home = "D:/GoogleDrive/Documents/Lectures/2024.1st/2024KSBi_BIML/data4practice" ;
setwd(home) ;
library(Seurat) ; library(ggplot2) ;

# where CellRanger outputs
cellrangers = dir(paste0(getwd(),"/CellRangerOuts")) ;
cellrangers
# "ctl.1"  "ctl.2"  "ctl.3"  "luad.1" "luad.2" "luad.3"

# load each CellRanger output and merge as a seurat.object
for (i in 1:length(cellrangers)){
  data.i = Read10X(data.dir = paste0(getwd(),"/CellRangerOuts/",cellrangers[i])) ;

  colnames(data.i) = paste0(cellrangers[i],".",colnames(data.i)) ;
  obj.i = CreateSeuratObject(counts= data.i, project="lung_obj", min.cells=3, min.features=300) ;
  obj.i$orig.ident = cellrangers[i] ;
  obj.i[["percent.mt"]] = PercentageFeatureSet(obj.i, "^MT-") ;

  cat(paste0("i = ",i," | ",cellrangers[i],"\n")) ;
  if(i==1){luadobj = obj.i} else {luadobj = merge(luadobj, obj.i)}
} ;
save(luadobj, file=paste0(home,"luadobj.rda")) ;
```

36

- 18 -

# Data loading

```
> luadobj
An object of class Seurat
20776 features across 21165 samples within 1 assay
Active assay: RNA (20776 features, 0 variable features)
>
> head(luadobj@meta.data)
                       orig.ident nCount_RNA nFeature_RNA percent.mt
ctl.1.AAACCCAGTTATGACC      ctl.1       6342         1947   9.350363
ctl.1.AAACCCAGTTCGAGCC      ctl.1       2255         1046   5.986696
ctl.1.AAACGAACAAGGCGTA      ctl.1      31132         4264  10.741359
ctl.1.AAACGAACATCTTCGC      ctl.1       3025         1153   7.206612
ctl.1.AAACGAAGTGCGTTTA      ctl.1       2186         1211   6.953339
ctl.1.AAACGAAGTTGGGCCT      ctl.1       2677         1176  10.646246
>
>
> table(luadobj@meta.data$orig.ident)

 ctl.1  ctl.2  ctl.3 luad.1 luad.2 luad.3
  3565   4511   3656   3670   3091   2672
>
```

# Filtering out poor quality cells

```
luadobj = subset(luadobj, subset = nFeature_RNA > 200 & percent.mt < 20) ;

Idents(luadobj) = "orig.ident"

VlnPlot(luadobj, features = c("nFeature_RNA", "nCount_RNA", "percent.mt"), ncol = 3)
```

# Normalize expression matrix and
# identify top variable genes

```
> # normalize expression matrix and identify top 2000 varialbe genes
> luadobj <- NormalizeData(luadobj, normalization.method = "LogNormalize", scale.factor = 10000)
Performing log-normalization
0%   10   20   30   40   50   60   70   80   90   100%
[----|----|----|----|----|----|----|----|----|----|
**************************************************|
>
> luadobj <- FindVariableFeatures(luadobj, selection.method = "vst", nfeatures = 2000)
Calculating gene variances
0%   10   20   30   40   50   60   70   80   90   100%
[----|----|----|----|----|----|----|----|----|----|
**************************************************|
Calculating feature variances of standardized and clipped values
0%   10   20   30   40   50   60   70   80   90   100%
[----|----|----|----|----|----|----|----|----|----|
**************************************************|
```

# Normalization

```
# Identify the 10 most highly variable genes
top10 <- head(VariableFeatures(luadobj), 10)

# plot variable features with and without labels
variable_plot <- VariableFeaturePlot(luadobj)
variable_plot.wTop10 <- LabelPoints(plot = variable_plot, points = top10, repel = TRUE)
variable_plot.wTop10
```

# Normalization and PCA

```
luadobj = ScaleData(luadobj, features=rownames(luadobj)) ;
## in case of removing unwanted sources of variation like MT contamination or cell cycling,
## regress out such heterogeneity
# luadobj = ScaleData(luadobj, features=rownames(luadobj), vars.to.regress="percent.mt") ;

luadobj = RunPCA(luadobj, features=VariableFeatures(object=luadobj), npcs=100) ;
DimPlot(luadobj, reduction = "pca")
plot(luadobj@reductions$pca@stdev)
```



# Estimate statistical significances of PCs

```
luadobj = JackStraw(luadobj, num.replicate=100, dims=100) ;
luadobj = ScoreJackStraw(luadobj, dims=1:100) ;
JackStrawPlot(luadobj, dims=1:100)
```

# Dimensional reduction and visualization

```r
luadobj = RunTSNE(luadobj, reduction="pca", dims=1:60) ;
luadobj = RunUMAP(luadobj, reduction="pca", dims=1:60) ;
DimPlot(luadobj, reduction = "tsne") ;
DimPlot(luadobj, reduction = "umap") ;
```



43

# Batch-correction



44

# Batch-correction

## Fast, sensitive and accurate integration of single-cell data with Harmony

Ilya Korsunsky [1,2,3,4], Nghia Millard[1,2,3,4], Jean Fan [5], Kamil Slowikowski[1,2,3,4], Fan Zhang [1,2,3,4], Kevin Wei[2], Yuriy Baglaenko [1,2,3,4], Michael Brenner[2], Po-ru Loh [1,3,4] and Soumya Raychaudhuri [1,2,3,4,6*]

---

# Batch-correction

```
> head(luadobj@meta.data)
                        orig.ident nCount_RNA nFeature_RNA percent.mt
ctl.1.AAACCCAGTTATGACC      ctl.1       6342        1947    9.350363
ctl.1.AAACCCAGTTCGAGCC      ctl.1       2255        1046    5.986696
ctl.1.AAACGAACAAGGCGTA      ctl.1      31132        4264   10.741359
ctl.1.AAACGAACATCTTCGC      ctl.1       3025        1153    7.206612
ctl.1.AAACGAAGTGCGTTTA      ctl.1       2186        1211    6.953339
ctl.1.AAACGAAGTTGGGCCT      ctl.1       2677        1176   10.646246
> table(luadobj@meta.data$orig.ident)

 ctl.1  ctl.2  ctl.3 luad.1 luad.2 luad.3
  3565   4511   3656   3670   3091   2672
> unique(luadobj@meta.data$orig.ident)
[1] "ctl.1"  "ctl.2"  "ctl.3"  "luad.1" "luad.2" "luad.3"
```

# Batch-correction

```r
library(harmony) ;
luadobj = RunHarmony(luadobj, group.by.vars="orig.ident") ;
ElbowPlot(luadobj, reduction="harmony", ndims=100) ;

luadobj = RunTSNE(luadobj, reduction="harmony", dims=1:60, seed.use=1234) ;
luadobj = RunUMAP(luadobj, reduction="harmony", dims=1:60, seed.use=1234) ;
DimPlot(luadobj, reduction = "tsne") ;
DimPlot(luadobj, reduction = "umap") ;

save(luadobj, file="luadobj.rda")
```

---

# Batch-correction



Batch-correction

https://drive.google.com/drive/folders/1R-5vQUaDBk59n-Iu5f635ANcHhBqw57j?usp=sharing

··· > 2024KSBi_BIML > data4practice ▾ 🔏

[ 유형 ▾ ] [ 사람 ▾ ] [ 수정 날짜 ▾ ]

| 이름 | 소유자 | 마지막으로 수정... ▾ ↓ | 파일 크기 |
|---|---|---|---|
| 📄 20240207.KSBi_BIML.practice.R 👥 | 🌑 나 | 오후 7:39 나 | 3KB |
| 📄 luadobj.rda 👥 | 🌑 나 | 오후 7:31 나 | 4.43GB |

# Thank you!

KIMQTAE@ajou.ac.kr

# KSBi-BIML 2024

Single-cell RNA-sequencing analysis:
Assignment of cell types (part2)

Kyu-Tae Kim
Ajou University School of Medicine

---

## 본 교육의 목표와 특징

### 단일세포 전사체 데이터 세포 종류 결정하기

- 클러스터링 분석의 의미를 이해한다.

- 클러스터링 종류와 방법을 이해한다.

- 세포 타입 결정 과정을 이해한다.

- 단일세포 전사체 데이터 clustering 과정을 이해한다.

- 단일세포 전사체로부터 cell type assignment 과정을 이해한다.

# How to understand thousands of individual things?



How to describe??

Given that we have individual pieces of fruits (single-cell analysis),
then how to sort these with which criteria? Color? Freshness? Kinds?

3

# Clustering objects



Supervised - clustering          by colors          by kinds          by freshness

[Red] cluster     [Green] cluster     [Grape] cluster     [Apple] cluster     [Rotten] cluster     [Fresh] cluster

4

# Supervised vs. Un-supervised clustering

| Supervised clustering | Un-supervised clustering |
|---|---|
| > The classes are predefined, and the task is to understand the basis for the classification from a set of labeled objects (training or learning set).<br><br>> This information is then used to classify future observations.<br><br>• Discriminant analysis<br>• Class prediction<br>• Supervised pattern recognition | > The classes are unknown a priori and need to be "discovered" from the data.<br><br><br><br>• Cluster analysis<br>• Class discovery<br>• Unsupervised pattern recognition |

5

# Clustering analysis

> **Finding groups** of objects such that the objects in a group will be similar (or related) to one another and different from (or unrelated to) the objects in other groups

> **Cluster**: a collection of 'similar' data



6

# Evaluation of clustering



> <u>Intra</u>-cluster distance:

the distance among members of a cluster

> <u>Inter</u>-cluster distance:

the distance between two different clusters

> A **good clustering** method will produce high quality clusters with

Low intra-class distance  =  High intra-class similarity

High inter-class distance  =  Low inter-class similarity

> How to determine **'similarity'**?

> How to measure **'distance'**?

---

# 클러스터링 종류와 방법

# Similarity measures with a gene expression table



**Data matrix**

> $n$: size of the data (how many samples)

> $p$: attributes of the data (how many genes)

$$\begin{bmatrix} x_{11} & \cdots & x_{1f} & \cdots & x_{1p} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ x_{i1} & \cdots & x_{if} & \cdots & x_{ip} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ x_{n1} & \cdots & x_{nf} & \cdots & x_{np} \end{bmatrix}$$

**Dimensionality**

$= n \times p$

**Distance matrix**
(dissimilarity matrix)

$$\begin{bmatrix} 0 \\ d(2,1) & 0 \\ d(3,1) & d(3,2) & 0 \\ \vdots & \vdots & \vdots \\ d(n,1) & d(n,1) & \cdots & \cdots & 0 \end{bmatrix}$$

---

# Classical distance measures

# Measures of relative distances

> **Pearson correlation**
- Measuring the degree of a linear relationship between two profiles

$$d_{cor}(x,y) = 1 - \frac{\sum\limits_{i=1}^{n}(x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum\limits_{i=1}^{n}(x_i - \bar{x})^2 \sum\limits_{i=1}^{n}(y_i - \bar{y})^2}}$$

**Parametric**

> **Eisen cosine correlation**
- A special case of Pearson's correlation with $x$ and $y$ both replaced by zero

$$d_{eisen}(x,y) = 1 - \frac{\left|\sum\limits_{i=1}^{n} x_i y_i\right|}{\sqrt{\sum\limits_{i=1}^{n} x_i^2 \sum\limits_{i=1}^{n} y_i^2}}$$

> **Spearman correlation**
- Measuring the correlation between the rank of $x$ and the rank of $y$ variables

$$d_{spear}(x,y) = 1 - \frac{\sum\limits_{i=1}^{n}(x_i' - \bar{x'})(y_i' - \bar{y'})}{\sqrt{\sum\limits_{i=1}^{n}(x_i' - \bar{x'})^2 \sum\limits_{i=1}^{n}(y_i' - \bar{y'})^2}}$$

**non-Parametric**

> **Kendall correlation**
- Measuring the correspondence between the ranking of $x$ and $y$ variables

$$d_{kend}(x,y) = 1 - \frac{n_c - n_d}{\frac{1}{2}n(n-1)}$$

11

---

# Clustering methods



> **Hierarchical clustering**

> **Partitioning clustering**
- K-medoids
- PAM (Partitioning Around Medoid)
- SOM (Self Organizing Maps)

> **Advanced clustering**
- Hybrid clustering methods
- Fuzzy clustering
- Model-based clustering
- Density-based clustering
- Graph-based clustering
- and ...

12

# Hierarchical clustering

**non-Hierarchical**



(Nested clusters)          (Dendrogram)

**Hierarchical**

---

# Hierarchical clustering

- Hierarchical clustering was the first algorithm used in microarray research to cluster genes. (David Bostein group, PNAS 1998)
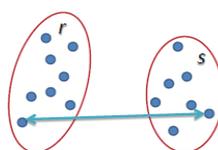


Average distance

- First, each object is assigned to its own cluster. Then, iteratively, the two most similar clusters are joined, representing a new node of the clustering tree. The similarity matrix is updated. This process is repeated until only a single cluster remains. (agglomerative clustering)



$$L(r,s) = \min(D(x_{ri}, x_{sj}))$$

**> Single linkage**

- Smallest distance

$$L(r,s) = \max(D(x_{ri}, x_{sj}))$$

**> Complete linkage**

- Largest distance

$$L(r,s) = \frac{1}{n_r n_s} \sum_{i=1}^{n_r} \sum_{j=1}^{n_s} D(x_{ri}, x_{sj})$$

**> Average linkage**

- Average distance

# Partitioning clustering



> **Hierarchical**

- Clustering is hierarchical decomposition (i.e., multiple levels)
- It can not correct erroneous merges or splits

> **Partitioning**

- It find mutually exclusive clusters of spherical shape
- It may use mean or medoid to represent cluster center
- It may effective for small- to medium-size data sets

---

# *K-means* clustering

- Number of cluster, *K*, must be specified
- Each cluster is associated with an averaged point (centroid)
- Each point is assigned to the cluster with the closest centroid

- Basic algorithm

  1: Select *K* points as the initial centroids.

  2: **repeat**

  3:  From *K* clusters by assigning all points to the closest centroid.

  4:  Recompute the centroid of each cluster.

  5: **until** The centroids does not change



Before K-Means

After K-Means

# Limitation of *K-means* clustering

• Applicable only when mean is defined, then what about categorical data?

• Need to specify *K*, the number of clusters, in advance

• Unable to handle noisy data and outliers

• Not suitable to discover clusters with



Differing sizes        Differing densities        Non-convex shapes

**Overcoming limitations**

• Using many clusters (i.e., high *K*)

• Using K-medoids, instead of k-means which is sensitive to outliers

---

# Dimensional reduction for visualization

## > Projection methods
• PCA (Principal Component Analysis)
• t-SNE (t-distributed Stochastic Neighbor Embedding)
• UMAP (Uniform Manifold Approximation and Projection)



PCA        t-SNE        UMAP

# Graph-based clustering



KNN = 5

KNN = 10

- Louvain community detection is applied to a shared-nearest-neighbor graph connecting the cells and finds tightly connected communities in the graph

- Increasing the number of neighbors when constructing the cell–cell graph indirectly decreases the resolution of graph-based clustering.
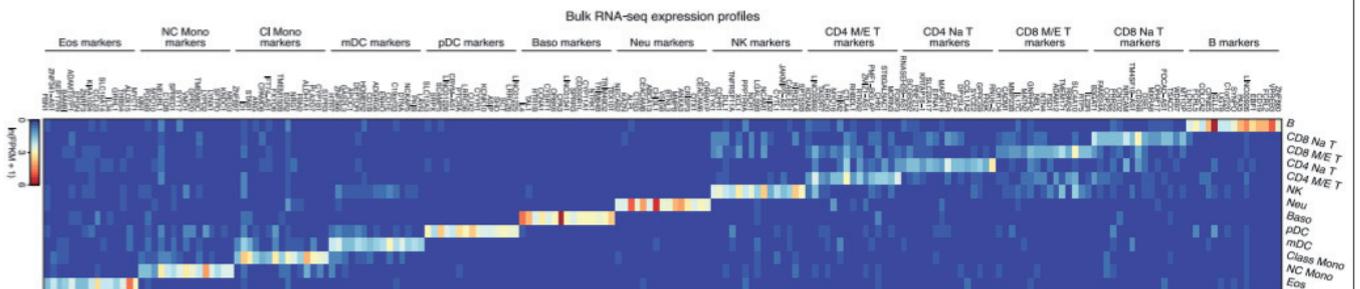
19

세 포 유 형 결 정

20

# Cell-type assignment



Litviňuková et al., *Nature* 2020

# Cell-type assignment



Litviňuková et al., *Nature* 2020

# Cell-type assignment



Travaglini et al., *Nature* 2020

23

# Cell-type assignment



Travaglini et al., *Nature* 2020

24

# Hematopoiesis

[cell type gene expression] https://dice-database.org/

[cell type gene expression] https://xteam.xbio.top/CellMarker/

[cell type gene expression] https://panglaodb.se/

[cell type gene expression] https://cellxgene.cziscience.com/cellguide/

[cell type gene expression] https://www.celltypist.org/encyclopedia/Immune/v2

[cell types in blood/tissue marker ptn. expression] https://www.proteinatlas.org/

# Useful resources to identify cell type markers

# Useful resources to identify cell type markers

# Data loading (practice)

https://drive.google.com/drive/folders/1R-5vQUaDBk59n-Iu5f635ANcHhBqw57j?usp=sharing

---

# Batch-correction
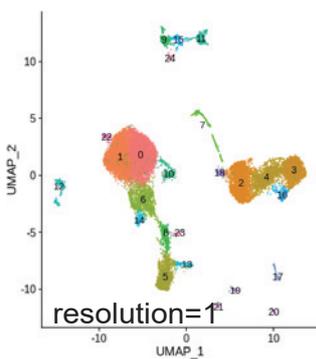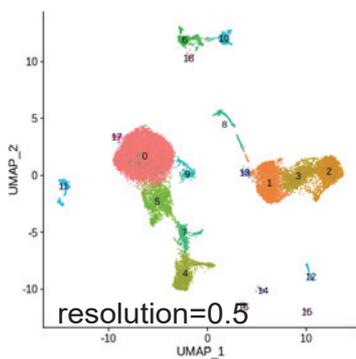


Batch-correction

# Clustering



```
> luadobj = FindNeighbors(luadobj, reduction="harmony", k.param=20, dims=1:60)
Computing nearest neighbor graph
Computing SNN
> luadobj = FindClusters(luadobj, resolution=1) ;
Modularity Optimizer version 1.3.0 by Ludo Waltman and Nees Jan van Eck

Number of nodes: 21165
Number of edges: 896920

Running Louvain algorithm...
0%   10   20   30   40   50   60   70   80   90   100%
[----|----|----|----|----|----|----|----|----|----|
*************************************************|
Maximum modularity in 10 random starts: 0.8588
Number of communities: 25
Elapsed time: 3 seconds
> DimPlot(luadobj, reduction="umap", group.by="seurat_clusters", pt.size=0.001, label=T)
```
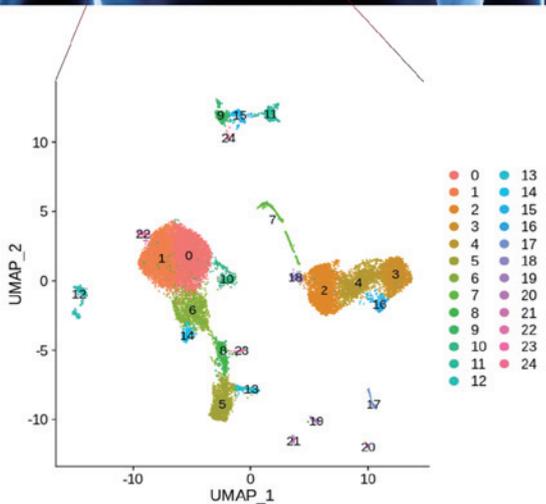


resolution=0.5

resolution=1

resolution=1.5

# Identification of cluster-specific markers

# Identification of cluster-specific markers

```
#MAST has good FDR control and is faster than DESeq2
luadobj.markers = FindAllMarkers(luadobj, only.pos=TRUE, min.pct=0.25, logfc.threshold=0.25, test.use="MAST") ;
```

test.use        Denotes which test to use. Available options are:

- "wilcox" : Identifies differentially expressed genes between two groups of cells using a Wilcoxon Rank Sum test (default)
- "bimod" : Likelihood-ratio test for single cell gene expression, (McDavid et al., Bioinformatics, 2013)
- "roc" : Identifies 'markers' of gene expression using ROC analysis. For each gene, evaluates (using AUC) a classifier built on that gene alone, to classify between two groups of cells. An AUC value of 1 means that expression values for this gene alone can perfectly classify the two groupings (i.e. Each of the cells in cells.1 exhibit a higher level than each of the cells in cells.2). An AUC value of 0 also means there is perfect classification, but in the other direction. A value of 0.5 implies that the gene has no predictive power to classify the two groups. Returns a 'predictive power' (abs(AUC-0.5) * 2) ranked matrix of putative differentially expressed genes.
- "t" : Identify differentially expressed genes between two groups of cells using the Student's t-test.
- "negbinom" : Identifies differentially expressed genes between two groups of cells using a negative binomial generalized linear model. Use only for UMI-based datasets
- "poisson" : Identifies differentially expressed genes between two groups of cells using a poisson generalized linear model. Use only for UMI-based datasets
- "LR" : Uses a logistic regression framework to determine differentially expressed genes. Constructs a logistic regression model predicting group membership based on each feature individually and compares this to a null model with a likelihood ratio test.
- "MAST" : Identifies differentially expressed genes between two groups of cells using a hurdle model tailored to scRNA-seq data. Utilizes the MAST package to run the DE testing.
- "DESeq2" : Identifies differentially expressed genes between two groups of cells based on a model using DESeq2 which uses a negative binomial distribution (Love et al, Genome Biology, 2014).This test does not support pre-filtering of genes based on average difference (or percent detection rate) between cell groups. However, genes may be pre-filtered based on their minimum detection rate (min.pct) across both cell groups. To use this method, please install DESeq2, using the instructions at https://bioconductor.org/packages/release/bioc/html/DESeq2.html

# Identification of cluster-specific markers

```
> head(luadobj.markers, 30)
            p_val avg_log2FC pct.1 pct.2 p_val_adj cluster       gene
INHBA           0   2.027373 0.957 0.324         0       0      INHBA
CCL20           0   1.948472 0.808 0.318         0       0      CCL20
CXCL3           0   1.835299 0.997 0.490         0       0      CXCL3
RND3            0   1.760417 0.711 0.210         0       0       RND3
TNF             0   1.717851 0.907 0.331         0       0        TNF
C1QA            0   1.593701 1.000 0.519         0       0       C1QA
IL1A            0   1.531175 0.692 0.180         0       0       IL1A
FBP1            0   1.530024 0.998 0.484         0       0       FBP1
FABP4           0   1.520869 0.967 0.406         0       0      FABP4
C1QB            0   1.489668 0.996 0.493         0       0       C1QB
CXCL5           0   1.463187 0.541 0.130         0       0      CXCL5
MCEMP1          0   1.349559 0.991 0.358         0       0     MCEMP1
SERPINA1        0   1.324133 0.998 0.501         0       0   SERPINA1
MRC1            0   1.294227 0.996 0.403         0       0       MRC1
ALDH2           0   1.290682 0.998 0.543         0       0      ALDH2
MARCO           0   1.281677 0.997 0.444         0       0      MARCO
SNX10           0   1.267388 0.989 0.432         0       0      SNX10
MS4A7           0   1.240686 0.998 0.449         0       0      MS4A7
VSIG4           0   1.233653 0.988 0.374         0       0      VSIG4
AC026369.3      0   1.210307 0.910 0.258         0       0 AC026369.3
LPL             0   1.186042 0.909 0.286         0       0        LPL
FTL             0   1.184815 1.000 0.996         0       0        FTL
C1QC            0   1.181961 0.989 0.374         0       0       C1QC
OLR1            0   1.164877 0.994 0.416         0       0       OLR1
STXBP2          0   1.144131 0.937 0.414         0       0     STXBP2
HLA-DRB5        0   1.140811 0.998 0.637         0       0   HLA-DRB5
LGALS3          0   1.119919 1.000 0.711         0       0     LGALS3
RETN            0   1.119641 0.794 0.301         0       0       RETN
MSR1            0   1.118809 0.983 0.391         0       0       MSR1
SERPING1        0   1.107077 0.944 0.352         0       0   SERPING1
```

# Identification of cluster-specific markers

```r
luadobj.markers.top20 = luadobj.markers %>% dplyr::group_by(cluster) %>% dplyr::top_n(n = 20, wt=avg_log2FC) ;

mySeuratClusters=unique(luadobj.markers.top20$cluster) ;

for(c in 1:length(mySeuratClusters)){
  luadobj.markers.top20.c = data.frame(
        cluster=luadobj.markers.top20[luadobj.markers.top20$cluster %in% mySeuratClusters[c], "gene"]) ;
  colnames(luadobj.markers.top20.c) = mySeuratClusters[c] ;
  if(c == 1){luadobj.markers.top20s = luadobj.markers.top20.c} else {
    luadobj.markers.top20s = cbind(luadobj.markers.top20s, luadobj.markers.top20.c)}
} ;
```
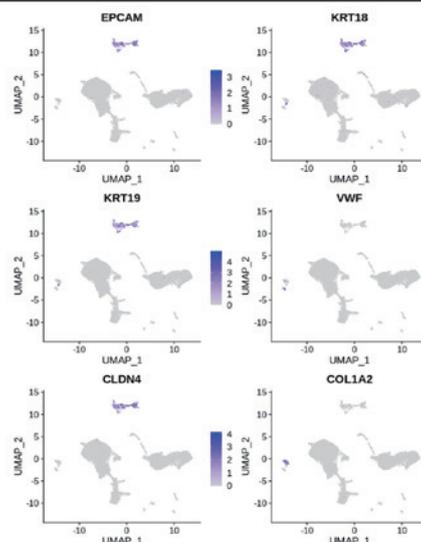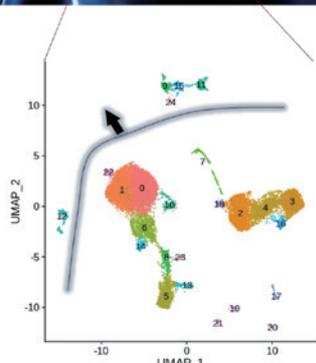
# Discovery of cluster identity



```r
nonImm.markers = c("EPCAM","KRT18","KRT19","VWF","CLDN4","COL1A2") ;
FeaturePlot(luadobj, features=nonImm.markers, reduction="umap") ;
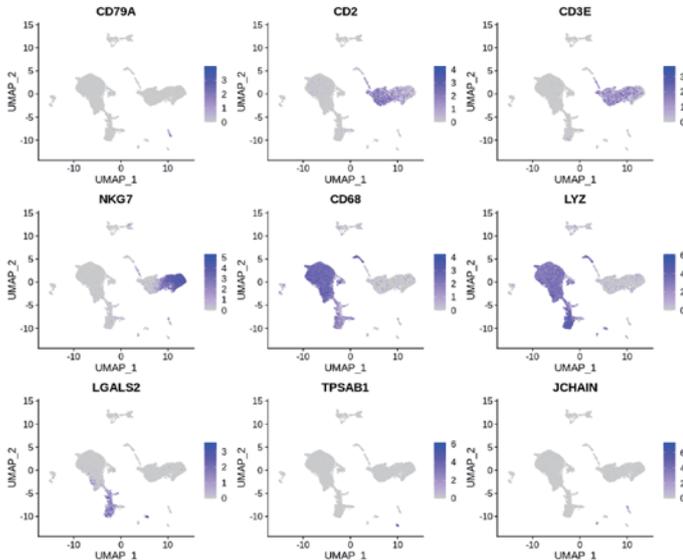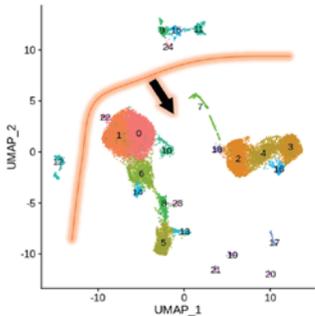```
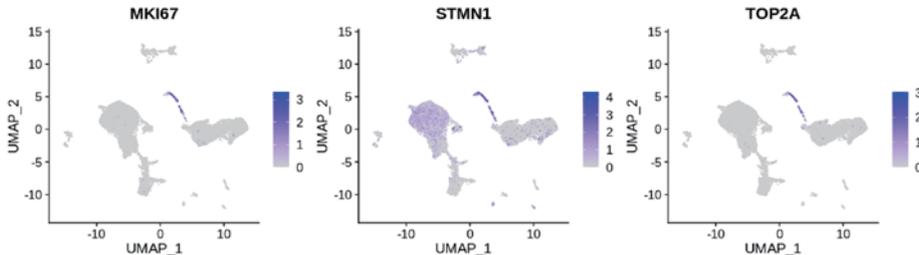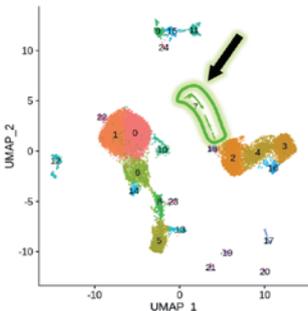
# Discovery of cluster identity



```
Imm.markers = c("CD79A","CD3E","NKG7","CD68","LYZ","LGALS2","TPSAB1","JCHAIN") ;
FeaturePlot(luadobj, features=Imm.markers, reduction="umap") ;
```

# Discovery of cluster identity



```
Proliferating.markers = c("MKI67","STMN1","TOP2A") ;
FeaturePlot(luadobj, features=Proliferating.markers, reduction="umap", ncol=3) ;
```
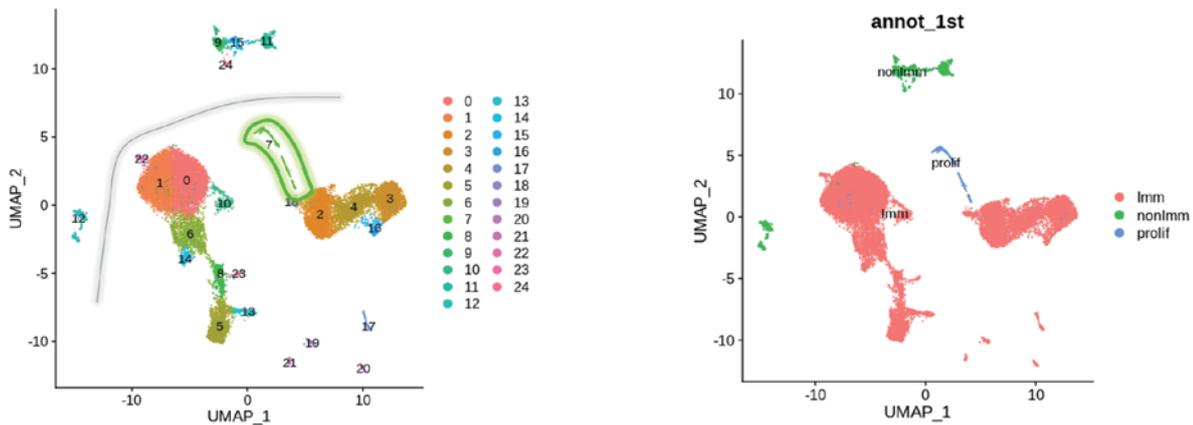
# Discovery of cluster identity

```
luadobj$annot_1st = "" ;
luadobj@meta.data[luadobj@meta.data$seurat_clusters %in% c(9,11,15,24, 12), ]$annot_1st <- "nonImm" ;
luadobj@meta.data[luadobj@meta.data$seurat_clusters %in% c(0:6,8,10,13,14,16:23), ]$annot_1st <- "Imm" ;
luadobj@meta.data[luadobj@meta.data$seurat_clusters %in% c(7), ]$annot_1st <- "prolif" ;


Idents(luadobj) = "annot_1st" ;
DimPlot(luadobj, group.by = "annot_1st", reduction="umap", label=T) ;
```
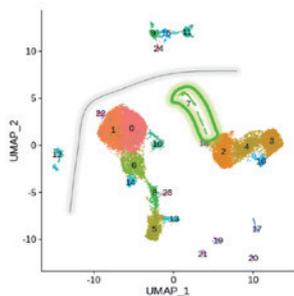
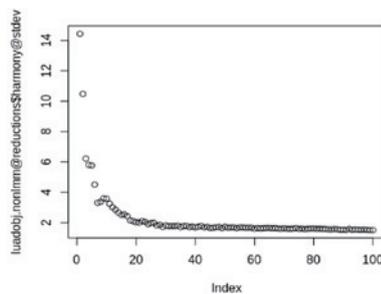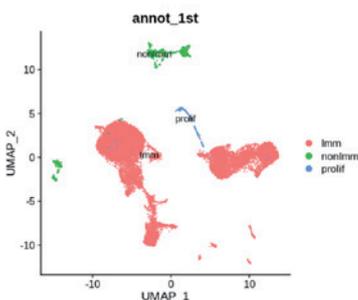# Discovery of sub-cluster identity



```
luadobj.nonImm = subset(luadobj, subset = seurat_clusters %in% c(9,11,15,24, 12)) ;

luadobj.nonImm = NormalizeData(luadobj.nonImm, normalization.method = "LogNormalize", scale.factor = 10000) ;
luadobj.nonImm = FindVariableFeatures(luadobj.nonImm, selection.method="vst", nfeatures=2000) ;
luadobj.nonImm = ScaleData(luadobj.nonImm, features=rownames(luadobj.nonImm)) ;
luadobj.nonImm = RunPCA(luadobj.nonImm, features=VariableFeatures(object=luadobj.nonImm), npcs=100) ;
luadobj.nonImm = RunHarmony(luadobj.nonImm, "orig.ident") ;
plot(luadobj.nonImm@reductions$harmony@stdev) ;
```
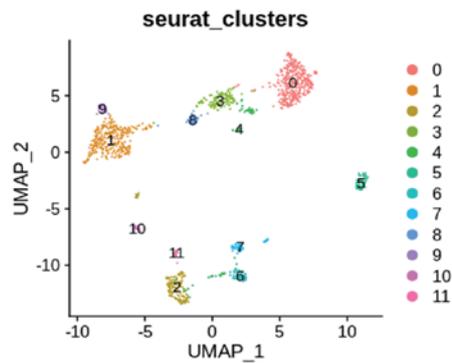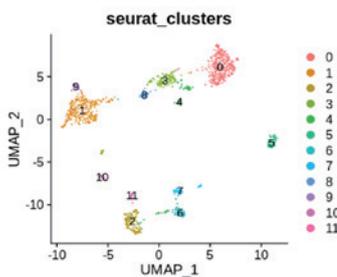
# Discovery of sub-cluster identity

```
luadobj.nonImm = RunUMAP(luadobj.nonImm, reduction="harmony", dims=1:40, seed.use=1234) ;
luadobj.nonImm = RunTSNE(luadobj.nonImm, reduction="harmony", dims=1:40, seed.use=1234) ;
luadobj.nonImm = FindNeighbors(luadobj.nonImm, reduction="harmony", dims=1:40)
luadobj.nonImm = FindClusters(luadobj.nonImm, resolution=0.5) ;
DimPlot(luadobj.nonImm, reduction="umap", group.by="seurat_clusters", pt.size=0.001, label=T) ;
```



41

# Discovery of sub-cluster identity



```
Epithelial.markers = c("EPCAM","CLDN4", "ELF3") ;
FeaturePlot(luadobj.nonImm, features=Epithelial.markers, reduction="umap", ncol=3) ;

Endothelial.markers = c("CDH5", "CLDN5", "RAMP2") ;
FeaturePlot(luadobj.nonImm, features=Endothelial.markers, reduction="umap", ncol=3) ;

Mesenchymal.markers = c("ITLN1","COL1A2","IGFBP6") ;
FeaturePlot(luadobj.nonImm, features=Mesenchymal.markers, reduction="umap", ncol=3) ;
```
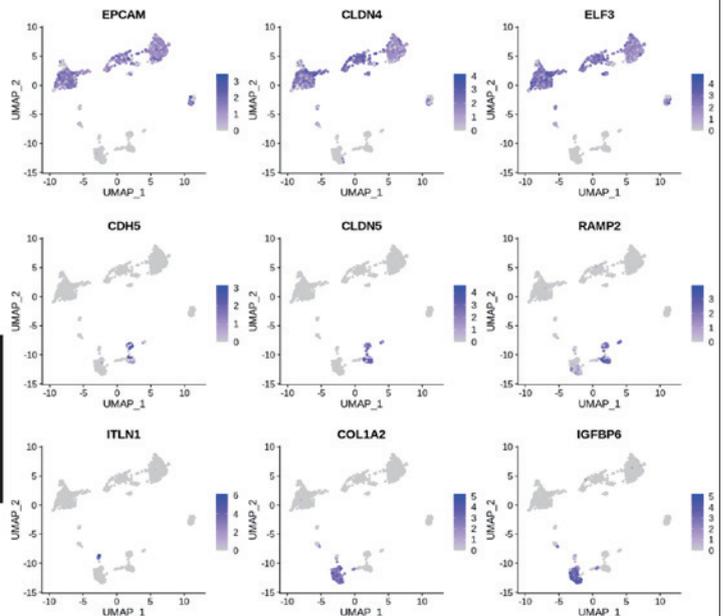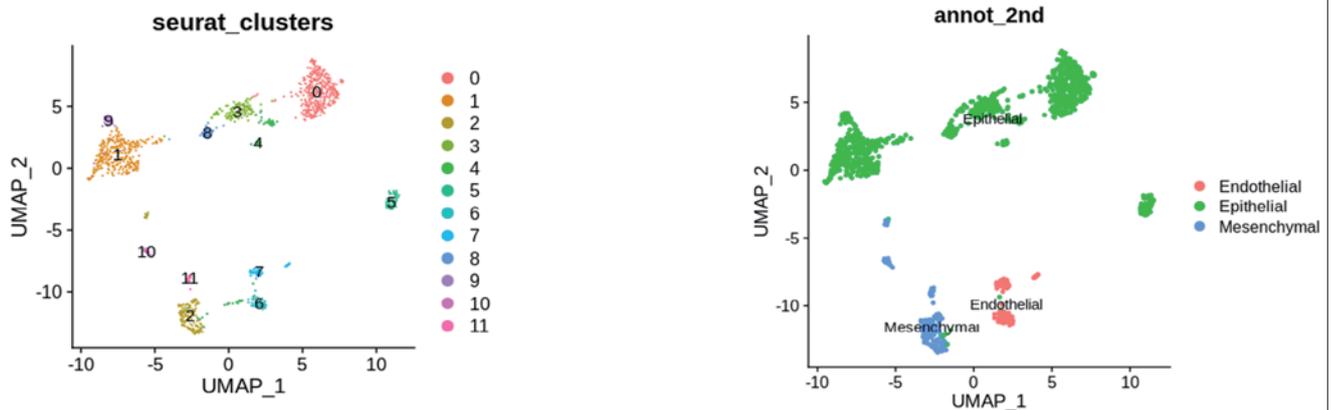
42

# Discovery of sub-cluster identity

```
luadobj.nonImm$annot_2nd = "" ;
luadobj.nonImm@meta.data[luadobj.nonImm@meta.data$seurat_clusters %in% c(1,9,8,3,4,0,5), ]$annot_2nd <- "Epithelial" ;
luadobj.nonImm@meta.data[luadobj.nonImm@meta.data$seurat_clusters %in% c(6,7), ]$annot_2nd <- "Endothelial" ;
luadobj.nonImm@meta.data[luadobj.nonImm@meta.data$seurat_clusters %in% c(2,10,11), ]$annot_2nd <- "Mesenchymal" ;

Idents(luadobj.nonImm) = "annot_2nd" ;
DimPlot(luadobj.nonImm, group.by = "annot_2nd", reduction="umap", label=T) ;
```
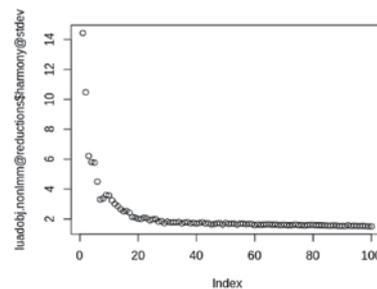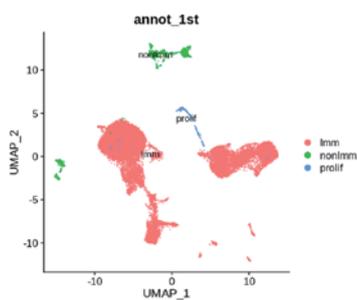


43

# Discovery of sub-cluster identity



```
luadobj.Imm = subset(luadobj, subset = seurat_clusters %in% c(0:6,8,10,13,14,16:23)) ;

luadobj.Imm = NormalizeData(luadobj.Imm, normalization.method = "LogNormalize", scale.factor = 10000) ;
luadobj.Imm = FindVariableFeatures(luadobj.Imm, selection.method="vst", nfeatures=2000) ;
luadobj.Imm = ScaleData(luadobj.Imm, features=rownames(luadobj.Imm)) ;
luadobj.Imm = RunPCA(luadobj.Imm, features=VariableFeatures(object=luadobj.Imm), npcs=100) ;
luadobj.Imm = RunHarmony(luadobj.Imm, "orig.ident") ;
plot(luadobj.Imm@reductions$harmony@stdev) ;
```

44

- 47 -

# Discovery of sub-cluster identity

```
luadobj.Imm = RunUMAP(luadobj.Imm, reduction="harmony", dims=1:50, seed.use=1234) ;
luadobj.Imm = RunTSNE(luadobj.Imm, reduction="harmony", dims=1:50, seed.use=1234) ;
luadobj.Imm = FindNeighbors(luadobj.Imm, reduction="harmony", dims=1:50)
luadobj.Imm = FindClusters(luadobj.Imm, resolution=0.8) ;
DimPlot(luadobj.Imm, reduction="umap", group.by="seurat_clusters", pt.size=0.001, label=T) ;
DimPlot(luadobj.Imm, reduction="tsne", group.by="seurat_clusters", pt.size=0.001, label=T) ;
```



45

# Discovery of sub-cluster identity

```
NK.markers = c("GNLY","KLRD1","KLRF1") ;
FeaturePlot(luadobj.Imm, features=NK.markers, reduction="umap", ncol=3) ;

T_common.markers = c("CD2","CD3D") ;
FeaturePlot(luadobj.Imm, features=T_common.markers, reduction="umap", ncol=3) ;

CD4.markers = c("CD4","CD40LG") ;
FeaturePlot(luadobj.Imm, features=CD4.markers, reduction="umap", ncol=3) ;

CD8.markers = c("CD8A","CD8B") ;
FeaturePlot(luadobj.Imm, features=CD8.markers, reduction="umap", ncol=3) ;

gdT.markers = c("TRDV2","TRGV9") ;
FeaturePlot(luadobj.Imm, features=gdT.markers, reduction="umap", ncol=3) ;

B.markers = c("CD79A","MS4A1","IGKC") ;
FeaturePlot(luadobj.Imm, features=B.markers, reduction="umap", ncol=3) ;

DC.markers = c("LGALS2","CPVL","CD1C") ;
FeaturePlot(luadobj.Imm, features=DC.markers, reduction="umap", ncol=3) ;

MQ.markers = c("MARCO","C1QA","FABP4") ;
FeaturePlot(luadobj.Imm, features=MQ.markers, reduction="umap", ncol=3) ;

Mono.markers = c("G0S2","S100A8","FCN1") ;
FeaturePlot(luadobj.Imm, features=Mono.markers, reduction="umap", ncol=3) ;

MAST.markers = c("TPSAB1","CPA3","GATA2") ;
FeaturePlot(luadobj.Imm, features=MAST.markers, reduction="umap", ncol=3) ;

pDC.markers = c("JCHAIN","IRF4","TCL1A") ;
FeaturePlot(luadobj.Imm, features=pDC.markers, reduction="umap", ncol=3) ;
```

46

# Discovery of sub-cluster identity

# Discovery of sub-cluster identity

```r
luadobj.Imm$annot_2nd = "" ;
luadobj.Imm@meta.data[luadobj.Imm@meta.data$seurat_clusters %in% c(2,12), ]$annot_2nd <- "NK" ;
luadobj.Imm@meta.data[luadobj.Imm@meta.data$seurat_clusters %in% c(3,8,13), ]$annot_2nd <- "CD4" ;
luadobj.Imm@meta.data[luadobj.Imm@meta.data$seurat_clusters %in% c(4), ]$annot_2nd <- "CD8" ;
luadobj.Imm@meta.data[luadobj.Imm@meta.data$seurat_clusters %in% c(14,21), ]$annot_2nd <- "B" ;
luadobj.Imm@meta.data[luadobj.Imm@meta.data$seurat_clusters %in% c(9,19,15), ]$annot_2nd <- "DC" ;
luadobj.Imm@meta.data[luadobj.Imm@meta.data$seurat_clusters %in% c(0,1,7,20,5,11,17), ]$annot_2nd <- "MQ" ;
luadobj.Imm@meta.data[luadobj.Imm@meta.data$seurat_clusters %in% c(6,10), ]$annot_2nd <- "Mono" ;
luadobj.Imm@meta.data[luadobj.Imm@meta.data$seurat_clusters %in% c(16), ]$annot_2nd <- "MAST" ;
luadobj.Imm@meta.data[luadobj.Imm@meta.data$seurat_clusters %in% c(18), ]$annot_2nd <- "pDC" ;

Idents(luadobj.Imm) = "annot_2nd" ;
DimPlot(luadobj.Imm, group.by = "annot_2nd", reduction="umap", label=T) ;
```
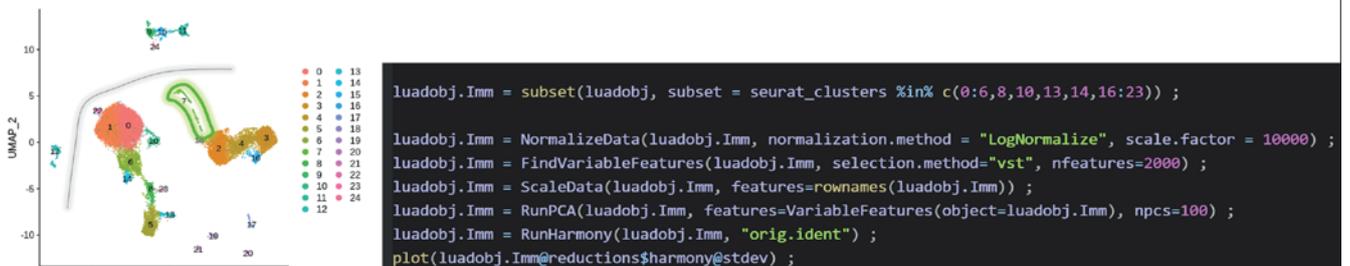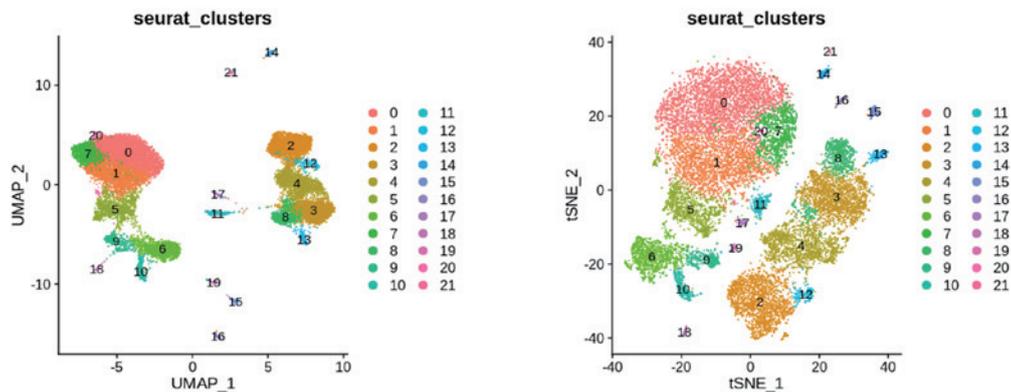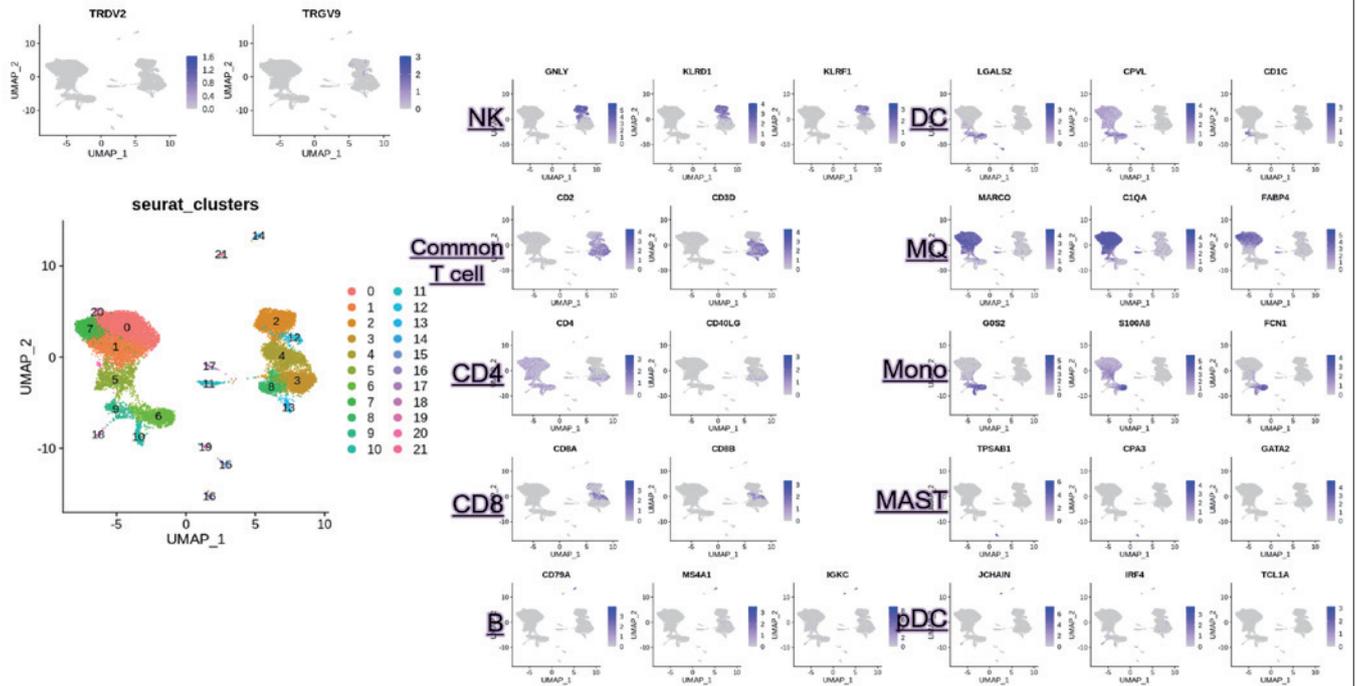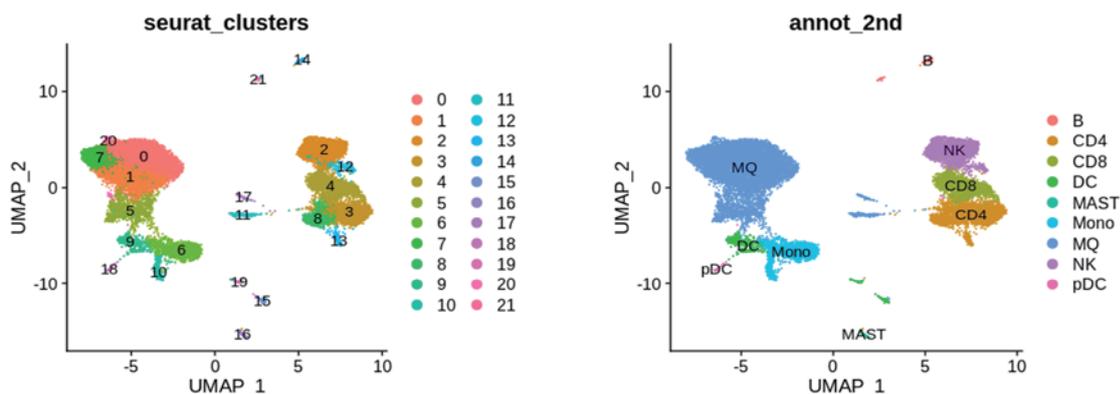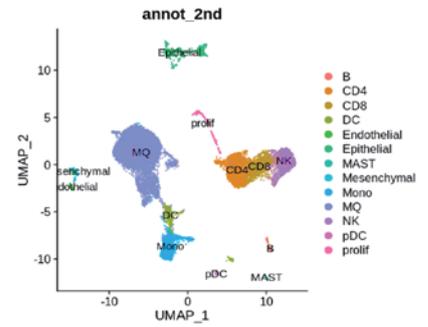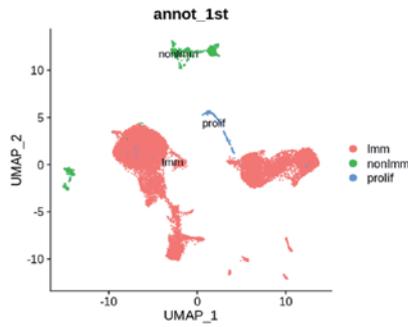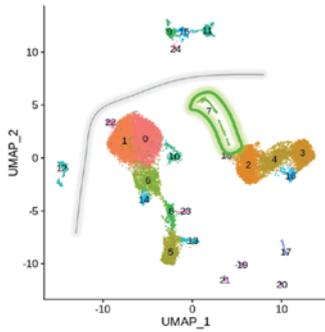
# Discovery of sub-cluster identity

```
luadobj$annot_2nd = "" ;
luadobj@meta.data[luadobj@meta.data$annot_1st %in% "prolif",]$annot_2nd <- "prolif" ;
luadobj@meta.data[rownames(luadobj@meta.data) %in% rownames(luadobj.nonImm@meta.data),]$annot_2nd <- luadobj.nonImm@meta.data$annot_2nd ;
luadobj@meta.data[rownames(luadobj@meta.data) %in% rownames(luadobj.Imm@meta.data),]$annot_2nd <- luadobj.Imm@meta.data$annot_2nd ;
```



49

---

# Thank you!

KIMQTAE@ajou.ac.kr