

# scRNAseq Analysis in R with Seurat

Monash Genomics and Bioinformatics Platform (MGBP)

Compiled: August 27, 2025



# Contents

<b>1 Getting started</b>	<b>5</b>
1.1 Summary . . . . .	5
<b>2 Set up</b>	<b>7</b>
2.1 Get the workshop material and data . . . . .	7
2.2 Package Installation . . . . .	7
2.3 Raw Data . . . . .	8
<b>3 Schedule</b>	<b>9</b>
<b>4 Load data</b>	<b>11</b>
4.1 Setup the Seurat Object . . . . .	11
4.2 The data set . . . . .	11
4.3 Different ways of loading the data . . . . .	12
<b>I Single Cell Analysis</b>	<b>15</b>
<b>5 QC Filtering</b>	<b>17</b>
5.1 QC and selecting cells for further analysis . . . . .	17
<b>6 Normalisation</b>	<b>27</b>
<b>7 PCAs and UMAPs</b>	<b>29</b>
7.1 Identification of highly variable features (feature selection) . . . . .	29
7.2 Scaling the data . . . . .	30
<b>8 Dimensionality reduction</b>	<b>33</b>
8.1 Perform linear dimensional reduction . . . . .	33
8.2 Determine the ‘dimensionality’ of the dataset . . . . .	38
8.3 Run non-linear dimensional reduction (UMAP/tSNE) . . . . .	39
8.4 Save . . . . .	40
<b>9 Data set integration with Harmony</b>	<b>43</b>

<b>10 Clustering</b>	<b>51</b>
10.1 Cluster cells . . . . .	51
10.2 Choosing a cluster resolution . . . . .	53
<b>11 Cluster Markers</b>	<b>57</b>
11.1 Finding differentially expressed features (cluster biomarkers) . . .	57
11.2 Use makers to label or find a cluster . . . . .	68
11.3 Assigning cell type identity to clusters . . . . .	69
<b>II Futher Analysis</b>	<b>73</b>
<b>12 SingleR</b>	<b>75</b>
<b>13 Differential Expression</b>	<b>79</b>
13.1 Prefiltering . . . . .	83
13.2 Default Wilcox test . . . . .	86
13.3 Seurat Negative binomial . . . . .	88
13.4 Pseudobulk . . . . .	89
<b>14 Cell cycle Assignment</b>	<b>93</b>
<b>III Other Resources</b>	<b>99</b>
<b>15 Resources</b>	<b>101</b>
15.1 Help and fruther Resources . . . . .	102
15.2 Data . . . . .	103
15.3 Analysis Tools . . . . .	104
15.4 Preprocessing Tools . . . . .	105
<b>IV Seurat Object</b>	<b>107</b>
<b>16 Structure</b>	<b>109</b>
16.1 Load an existing Seurat object . . . . .	109
16.2 What's in there? . . . . .	111
<b>17 Acknowledgements</b>	<b>113</b>
<b>18 Session info</b>	<b>115</b>

---

# Chapter 1

## Getting started

Instructors: Paul Harrison & Laura Perlaza-Jimenez

### 1.1 Summary

This workshop, conducted by the Monash Genomics and Bioinformatics Platform, will cover how to extend analysis to contemporary third-party tools, Seurat, Harmony and SingleR. We will be walking through the Single Cell analysis using Seurat package, correction of batch effects with Harmony and cell annotation using SingleR.

If you are using virtual machine (VM) provided by the workshop, you don't need to install any R packages. **Skip SET UP** section.

If you are using your own computer, follow Installation and Setup instructions.

#### 1.1.1 Recommended Computer Requirements:

System Requirements:

*Windows:*

- Windows 8.1 (64-bit) or later
- 4GB RAM
- SSD storage highly recommended
- Updated video/display drivers recommended

*macOS:*

- macOS 10.15 (Catalina) or later
- 4GB RAM
- SSD storage highly recommended

Install latest versions of:

- R
- RStudio
- Seurat
- SeuratWrappers
- SingleR
- Celldex
- Harmony

# Chapter 2

## Set up

*For this workshop you will have available a VM.* However, if you which to use your local computer here are the set up instructions:

**IMPORTANT:** If you have an M1 Mac - Make sure you have gfortan.

### 2.1 Get the workshop material and data

In RStudio **create a new project**. This ensures all the files for this workshop are placed in their own folder.

Once you've created a new project, run the following R code

### 2.2 Package Installation

For this workshop, several packages need to be installed.

BiocManager likes to update installed packages, but we have disabled this in the R code below. If your installation fails, then you might need to turn updates on. Note that if there are a large number of packages that BiocManager wants to update it can take several hours.

These instructions have been tested with R version 4.4.0 and Bioconductor version 3.16.

```
## Install required packages for Seurat and clustree:  
install.packages(c("Seurat", "dplyr", "remotes", "R.utils", "harmony",  
    "hdf5r", "clustree", "RColorBrewer", "tidyverse", "pander"))  
  
## Install required Bioconductor packages  
install.packages("BiocManager")  
BiocManager::install(c('SingleR', 'celldex',
```

```
'BiocGenerics', 'DelayedArray', 'DelayedMatrixStats',
'limma', 'S4Vectors', 'SingleCellExperiment',
'SummarizedExperiment','edgeR' ),
update=FALSE)

devtools::install_github('immunogenomics/presto@1.0.0')
```

## 2.3 Raw Data

```
# where are you? what folder are you working in

getwd()
#> [1] "/mnt/ceph/mpb/servers/bioinformatics-platform/home/lper0012/tasks/training/scR

## Download and untar the data
options(timeout=3600)
download.file(
  "https://bioinformatics.erc.monash.edu/home/lper0012/SingleCellWorkshopData/single"

untar("data.tar")

## Download workshop.R script. This file contains all the example code for this workshop
download.file(
  "https://monashbioinformaticsplatform.github.io/scRNAseq_Workshop_ABACBS_2024/work
```

# Chapter 3

## Schedule

### Day 1

Time	Content
10:00	Welcome and Housekeeping Log in to VM Background on ScRNAseq + R tools Load data
11:00-11:15	Break (15 mins)
11:15	QC filtering Normalisation
12:15 -12:45	Lunch Break PCAs and UMAPs
14:00	Finish Day 1

### Day 2

Time	Content
10:00	Integration
11:00-11:15	Break (15 mins)
11:15	Clustering Cluster Markers SingleR
12:15-12:45	Break (15mins) Differential Expression Main room Discussion on differential expression
14:00	Final remarks and conclusion



# Chapter 4

## Load data

Let's get started with a single cell introduction

### 4.1 Setup the Seurat Object

### 4.2 The data set

The dataset used in this workshop is a *modified* version derived from this study (see here). It has been adapted to introduce additional complexity for instructional purposes. *Please refrain from drawing any biological conclusions from this data as it does not represent real experimental results.*

This dataset represents human peripheral blood mononuclear cells (PBMCs), pooled from eight individual donors. Genetic differences among donors enable the identification of some cell doublets, enhancing data complexity. It includes two single-cell sequencing batches, one of which was stimulated with IFN-beta. Additionally, mitochondrial expression levels have been introduced to demonstrate how to interpret and apply mitochondrial thresholds for filtering purposes.

#### Note: What does the data look like?

What do the input files look like? It varies, but this is the output of the Cell-Ranger pipeline, described here

```
analysis
  clustering
  diffexp
  pca
  tsne
  umap
cloupe.cloupe
```

```

filtered_feature_bc_matrix
  barcodes.tsv.gz
  features.tsv.gz
  matrix.mtx.gz
filtered_feature_bc_matrix.h5. --> matrix to read in h5 format
metrics_summary.csv
molecule_info.h5
possorted_genome_bam.bam
possorted_genome_bam.bam.bai
raw_feature_bc_matrix
  barcodes.tsv.gz
  features.tsv.gz
  matrix.mtx.gz
raw_feature_bc_matrix.h5
web_summary.html

```

We start by reading in the data. There are several options for loading the data. The `Read10X()` function reads in the output of the cellranger pipeline from 10X, returning a unique molecular identified (UMI) count matrix. The values in this matrix represent the number of molecules for each feature (i.e. gene; row) that are detected in each cell (column).

We next use the count matrix to create a `Seurat` object. The object serves as a container that contains both data (like the count matrix) and analysis (like PCA, or clustering results) for a single-cell dataset. For a technical discussion of the `Seurat` object structure, check out the GitHub Wiki. For example, the count matrix is stored in `pbmc@assays$RNA@counts`.

```

library(dplyr)
library(ggplot2)
library(Seurat)
library(patchwork)

```

### 4.3 Different ways of loading the data

Example 1. Load your data using the path to the folder: `filtered_feature_bc_matrix` that is in the output folder of your cellranger run using the `Read10X` function.

```

# Load the PBMC dataset
pbmc.data <- Read10X(data.dir = "outs/filtered_feature_bc_matrix")
# Initialize the Seurat object with the raw (non-normalized data).
seurat_object <- CreateSeuratObject(counts = pbmc.data, min.cells = 3, min.features =

```

Example 2. Load your data directing the `ReadMtx` function to each of the relevant files in the `filtered_feature_bc_matrix` folder in the outputs from

your cellranger run. MTX is a simple text format for sparse matrices.

```
expression_matrix <- ReadMtx(
  mtx = "outs/filtered_feature_bc_matrix/count_matrix.mtx.gz", features = "outs/filtered_feature_bc_matrix/barcodes.tsv.gz"
)
seurat_object <- CreateSeuratObject(counts = expression_matrix)
```

Example 3. Load your data using the `Read10X_h5` function to each of the relevant HDF5 files. HDF5 is an efficient binary format.

```
pbmc.data <- Read10X_h5("data/filtered_feature_bc_matrix.h5")
metadata <- read.table("data/metadata.txt")

seurat_object <- CreateSeuratObject(counts = pbmc.data ,
                                      assay = "RNA", project = 'pbmc')
#> Warning: Feature names cannot have underscores ('_'),
#> replacing with dashes ('-')

# Use AddMetaData to add new meta data to object
seurat_object <- AddMetaData(object = seurat_object, metadata = metadata)
```

What does data in a count matrix look like?

```
# Lets examine a few genes in the first thirty cells
pbmc.data[c("CD3D", "TCL1A", "MS4A1"), 1:30]
#> 3 x 30 sparse Matrix of class "dgCMatrix"
#> [[ suppressing 30 column names 'AGGGCGCTATTC-1', 'GGAGACGATTGTT-1', 'CACCGTTGTCGTAG-1' ...
#>
#> CD3D . 2 . . . 5 . . . 2 . . . 1 . . 1 1 . . . 1 . 1 3 .
#> TCL1A . . . . . . . . . . . . . . . . . . . . . . . . . 3 . .
#> MS4A1 . . . . . . . . . . . . . . . . . . . . . . . . . 1 . .
#>
#> CD3D 2 . .
#> TCL1A . .
#> MS4A1 . . .
```

The . values in the matrix represent 0s (no molecules detected). Since most values in an scRNA-seq matrix are 0, Seurat uses a sparse-matrix representation whenever possible. This results in significant memory and speed savings for Drop-seq/inDrop/10x data.

```
dense.size <- object.size(as.matrix(pbmc.data))
#> Warning in asMethod(object): sparse->dense coercion:
#> allocating vector of size 1.3 GiB
dense.size
#> 1428252816 bytes
sparse.size <- object.size(pbmc.data)
```

```
sparse.size  
#> 39615768 bytes  
dense.size / sparse.size  
#> 36.1 bytes
```

# **Part I**

# **Single Cell Analysis**



# Chapter 5

## QC Filtering

The steps below encompass the standard pre-processing workflow for scRNA-seq data in Seurat. These represent the selection and filtration of cells based on QC metrics, data normalization and scaling, and the detection of highly variable features.

### 5.1 QC and selecting cells for further analysis

#### Why do we need to do this?

Low quality cells can add noise to your results leading you to the wrong biological conclusions. Using only good quality cells helps you to avoid this. Reduce noise in the data by filtering out low quality cells such as dying or stressed cells (high mitochondrial expression) and cells with few features that can reflect empty droplets.

Seurat allows you to easily explore QC metrics and filter cells based on any user-defined criteria. A few QC metrics commonly used by the community include

- The number of unique genes detected in each cell.
  - Low-quality cells or empty droplets will often have very few genes
  - Cell doublets or multiplets may exhibit an aberrantly high gene count
- Similarly, the total number of molecules detected within a cell (correlates strongly with unique genes)
- The percentage of reads that map to the mitochondrial genome
  - Low-quality / dying cells often exhibit extensive mitochondrial contamination
  - We calculate mitochondrial QC metrics with the `PercentageFeatureSet()` function, which calculates the percentage of counts originating from

a set of features

- We use the set of all genes starting with MT- as a set of mitochondrial genes

```
# The $ operator can add columns to object metadata.
# This is a great place to stash QC stats
seurat_object$percent.mt <- PercentageFeatureSet(seurat_object, pattern = "^\u00d7MT-")
```

### Challenge: The meta.data slot in the Seurat object

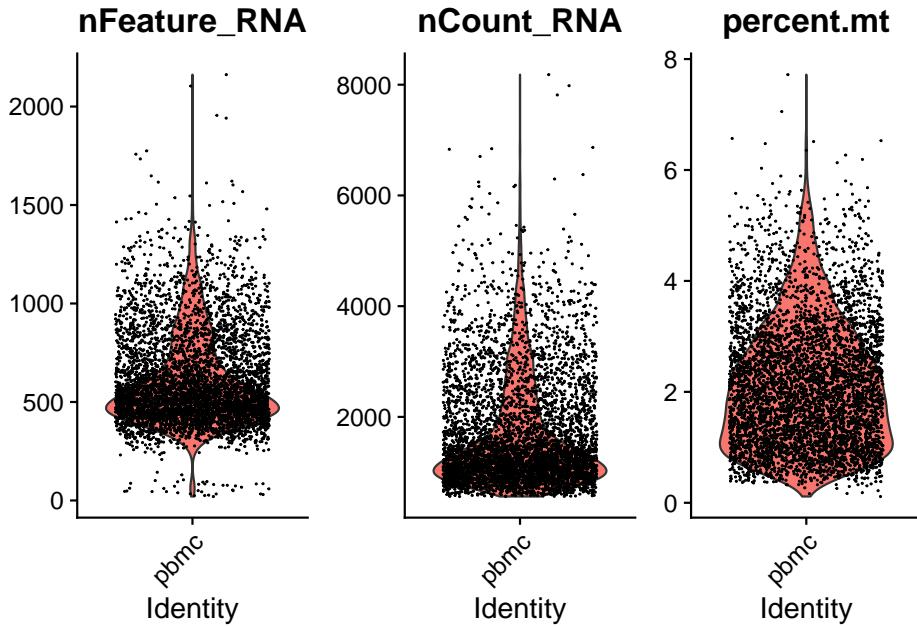
Where are QC metrics stored in Seurat?

- The number of unique genes and total molecules are automatically calculated during `CreateSeuratObject()`
  - You can find them stored in the object meta data

What do you notice has changed within the `meta.data` table now that we have calculated mitochondrial gene proportion?

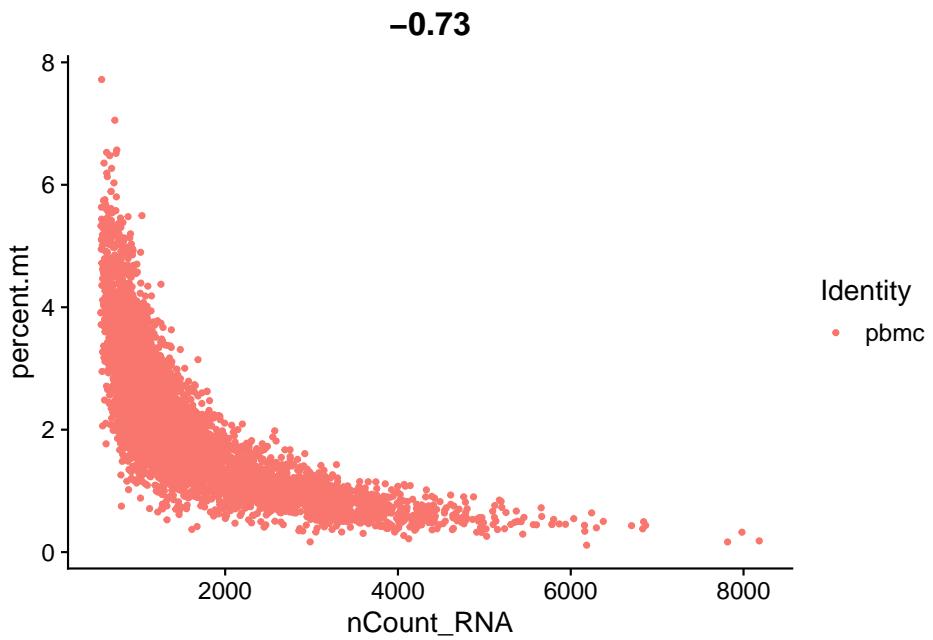
In the example below, we visualize QC metrics. We will use these to filter cells.

```
# Visualize QC metrics as a violin plot
VlnPlot(seurat_object, features = c("nFeature_RNA", "nCount_RNA", "percent.mt"), ncol =
#> Warning: Default search for "data" layer in "RNA" assay
#> yielded no results; utilizing "counts" layer instead.
#> Warning: The `slot` argument of `FetchData()` is deprecated as of
#> SeuratObject 5.0.0.
#> i Please use the `layer` argument instead.
#> i The deprecated feature was likely used in the Seurat
#> package.
#> Please report the issue at
#> <https://github.com/satijalab/seurat/issues>.
#> This warning is displayed once every 8 hours.
#> Call `lifecycle:::last_lifecycle_warnings()` to see where
#> this warning was generated.
#> Warning: `PackageCheck()` was deprecated in SeuratObject 5.0.0.
#> i Please use `rlang::check_installed()` instead.
#> i The deprecated feature was likely used in the Seurat
#> package.
#> Please report the issue at
#> <https://github.com/satijalab/seurat/issues>.
#> This warning is displayed once every 8 hours.
#> Call `lifecycle:::last_lifecycle_warnings()` to see where
#> this warning was generated.
```

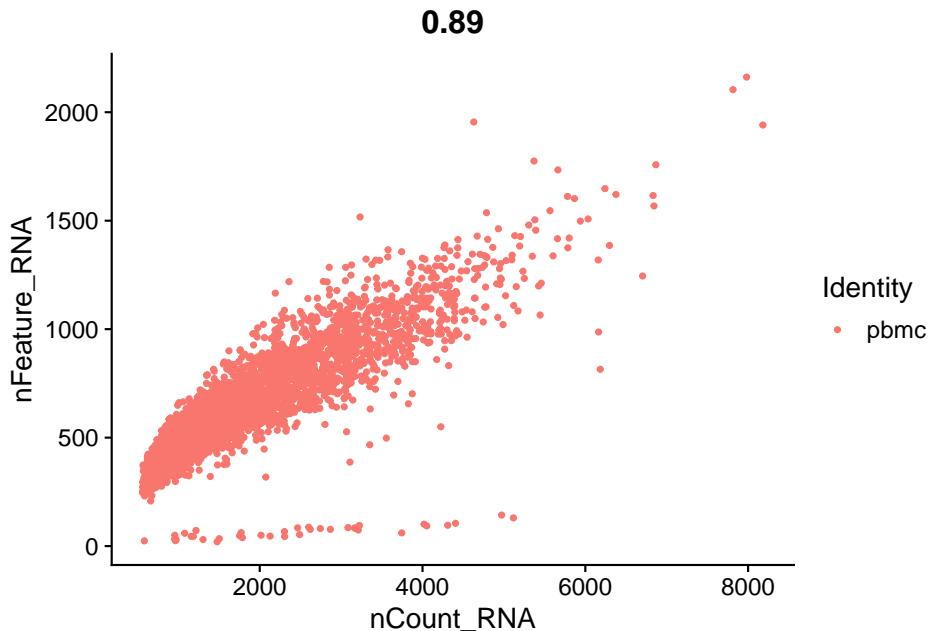


```
# FeatureScatter is typically used to visualize feature-feature relationships,
# but can be used for anything calculated by the object,
# i.e. columns in object metadata, PC scores etc.
```

```
FeatureScatter(seurat_object, feature1 = "nCount_RNA", feature2 = "percent.mt")
```



```
FeatureScatter(seurat_object, feature1 = "nCount_RNA", feature2 = "nFeature_RNA")
```



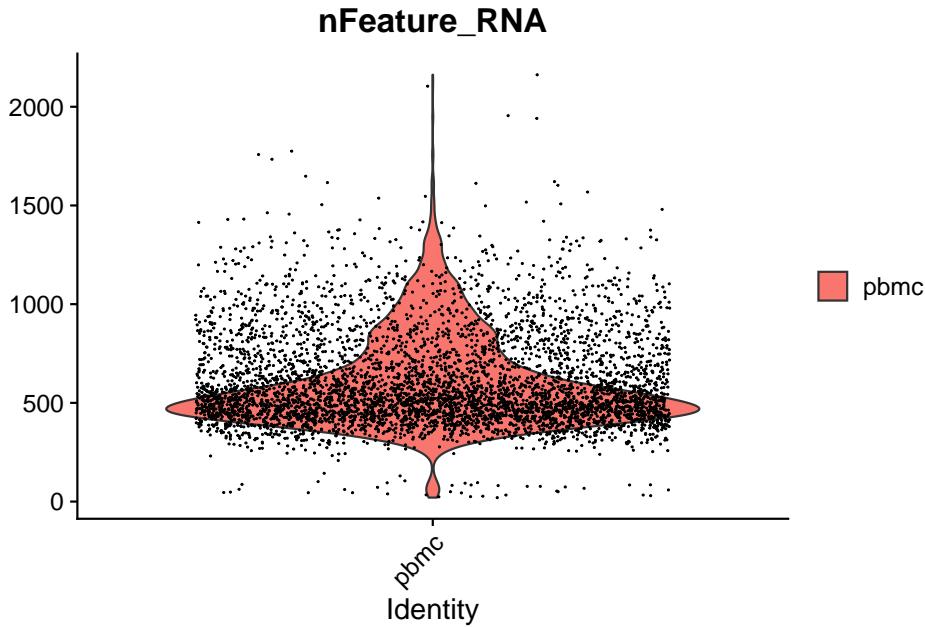
### Challenge: Red Blood Cells

Genes “HBA1”, “HBA2”, and “HBB” are components of hemoglobin in red blood cells.

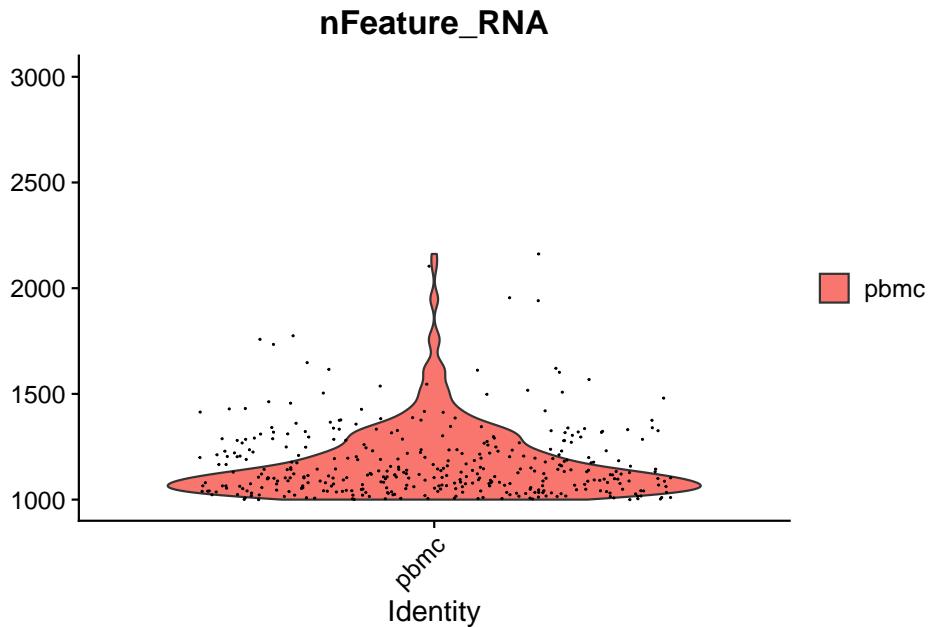
- Use PercentageFeatureSet, passing these genes to the “features” argument, to find cells that might be red blood cells.
- How do cells high in these genes differ from other cells, in terms of number of features or total count?
- Should we remove these cells?

Let's look again at the number of features (genes) to the percent mitochondrial genes plot.

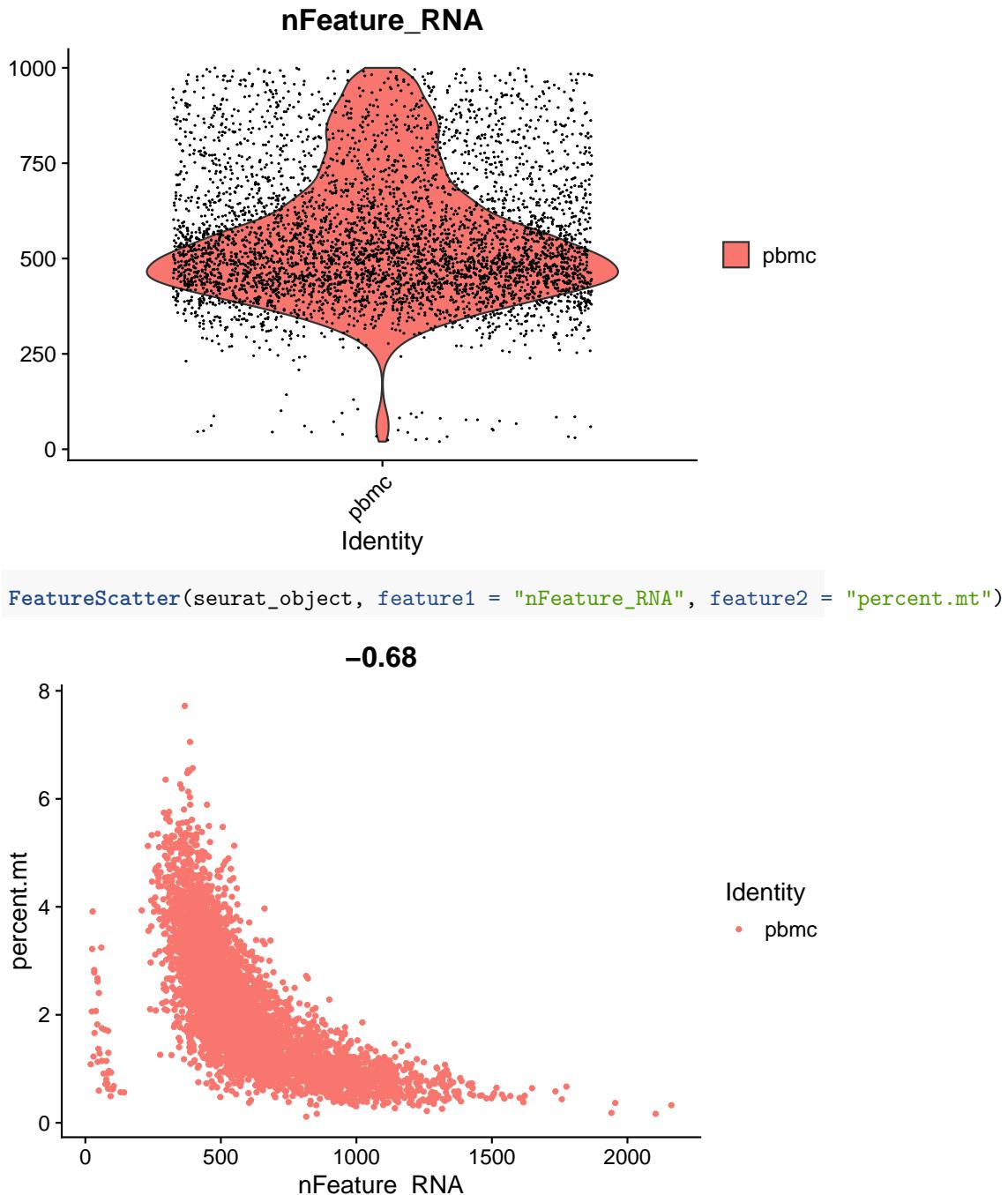
```
VlnPlot(seurat_object, features = "nFeature_RNA")
#> Warning: Default search for "data" layer in "RNA" assay
#> yielded no results; utilizing "counts" layer instead.
```



```
# Zoom in to the max and min
VlnPlot(seurat_object, features = "nFeature_RNA") + scale_y_continuous(limits = c(1000,3000))
#> Warning: Default search for "data" layer in "RNA" assay
#> yielded no results; utilizing "counts" layer instead.
#> Scale for y is already present.
#> Adding another scale for y, which will replace the existing
#> scale.
#> Warning: Removed 4624 rows containing non-finite outside the scale
#> range (`stat_ydensity()`).
#> Warning: Removed 4624 rows containing missing values or values
#> outside the scale range (`geom_point()`).
```



```
VlnPlot(seurat_object, features = "nFeature_RNA", y.max =1000)
#> Warning: Default search for "data" layer in "RNA" assay
#> yielded no results; utilizing "counts" layer instead.
#> Warning: Removed 376 rows containing non-finite outside the scale
#> range (`stat_ydensity()`).
#> Warning: Removed 376 rows containing missing values or values
#> outside the scale range (`geom_point()`).
```



You can check different thresholds of mito percentage.

```
#Number of cells that would be left after filters
# Proportion of cells with less than 5% mito

mean(seurat_object$percent.mt < 5)
#> [1] 0.983

# Proportion of cells with less than 2% mito

mean(seurat_object$percent.mt < 2)
#> [1] 0.5158
```

Ok, let's go with these filters:

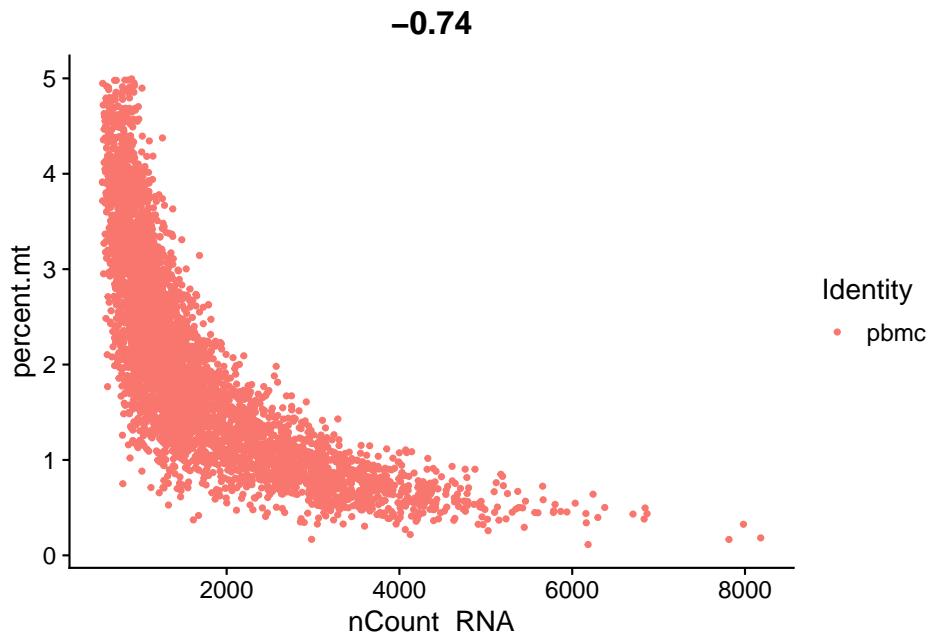
- We filter cells to have  $>200$  unique features
- We filter cells that have  $>5\%$  mitochondrial counts

Let's apply this and subset our data. This will remove the cells we think are of poor quality.

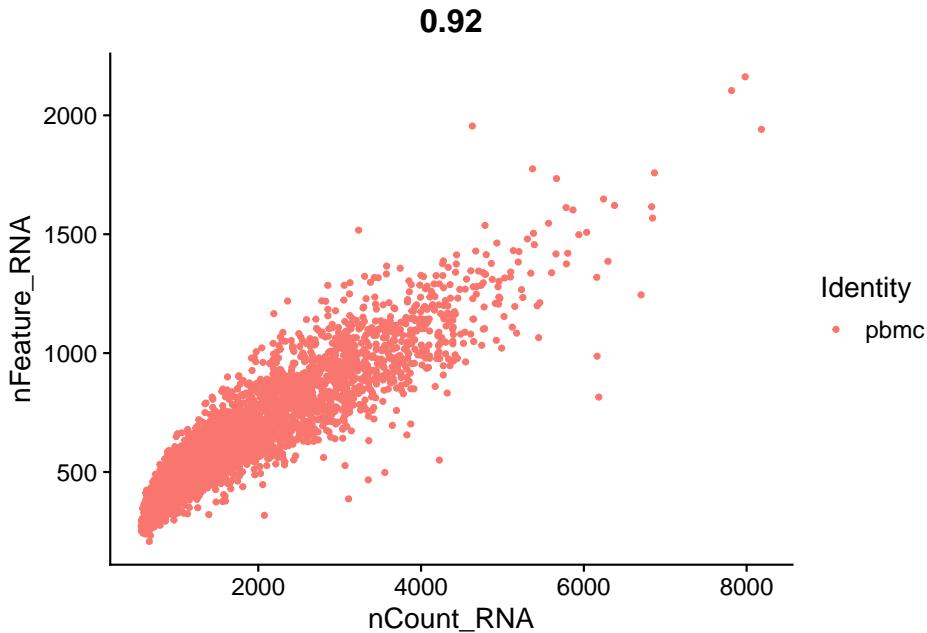
```
seurat_object <- subset(seurat_object, subset = nFeature_RNA > 200 & percent.mt < 5)
```

Let's replot the feature scatters and see what they look like.

```
FeatureScatter(seurat_object, feature1 = "nCount_RNA", feature2 = "percent.mt")
```

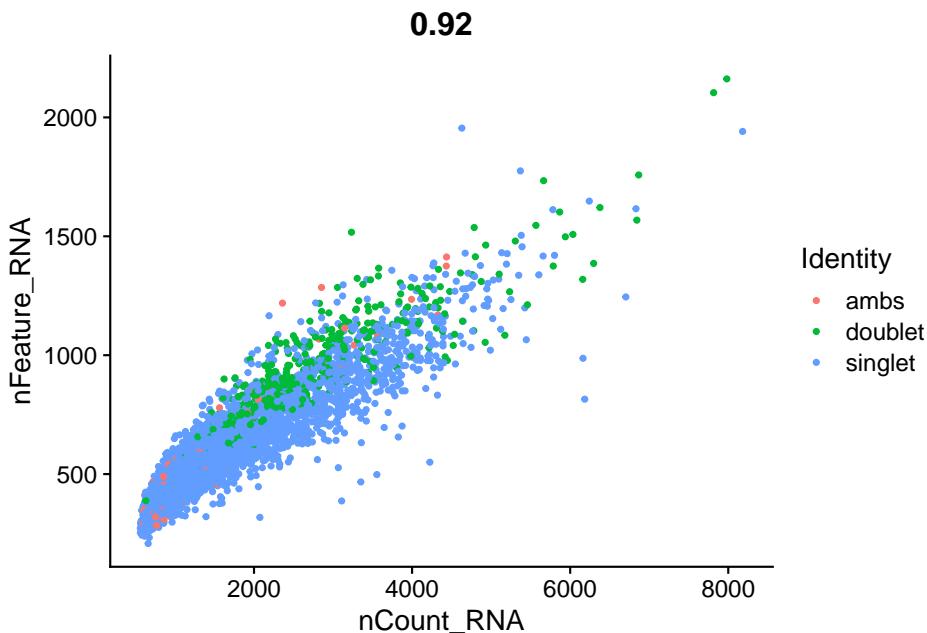


```
FeatureScatter(seurat_object, feature1 = "nCount_RNA", feature2 = "nFeature_RNA")
```



We also wondered if cells with high counts might be doublets. Should we also filter cells with very high counts? With this data, we know for certain some of the doublets!

```
FeatureScatter(seurat_object, feature1 = "nCount_RNA", feature2 = "nFeature_RNA", group.by="multi")
```





# Chapter 6

## Normalisation

slides

### Why do we need to do this?

The sequencing depth can be different per cell. This can bias the counts of expression showing higher numbers for more sequenced cells leading to the wrong biological conclusions. To correct this the feature counts are normalized.

After removing unwanted cells from the dataset, the next step is to normalize the data. By default, we employ a global-scaling normalization method “LogNormalize” that normalizes the feature expression measurements for each cell by the total expression, multiplies this by a scale factor (10,000 by default), and log-transforms the result. Normalized values are stored in `seurat_object$RNA@data`.

```
seurat_object <- NormalizeData(seurat_object, normalization.method = "LogNormalize", scale.factor  
#> Normalizing layer: counts
```

For clarity, in this previous line of code (and in future commands), we provide the default values for certain parameters in the function call. However, this isn’t required and the same behavior can be achieved with:

```
seurat_object <- NormalizeData(seurat_object)
```

There are other options for normalization such as SCTtransform which was popularized in 2019, however Log base normalization continued to be preferred as they perform better. See here for more details.



# Chapter 7

## PCAs and UMAPs

Slides

### 7.1 Identification of highly variable features (feature selection)

#### Why do we need to do this?

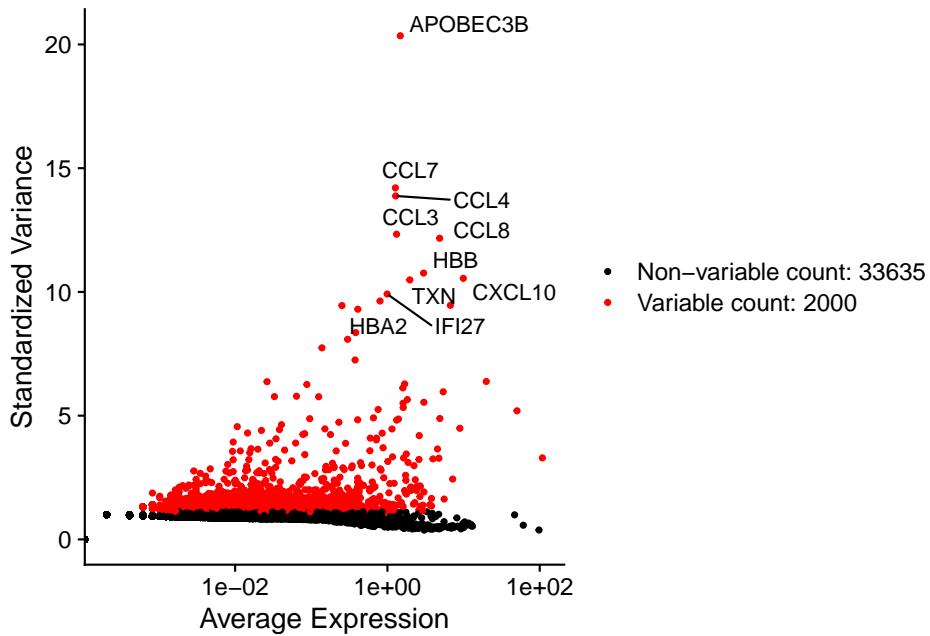
Identifying the most variable features allows retaining the real biological variability of the data and reduce noise in the data.

We next calculate a subset of features that exhibit high cell-to-cell variation in the dataset (i.e, they are highly expressed in some cells, and lowly expressed in others). The Seurat authors and others have found that focusing on these genes in downstream analysis helps to highlight biological signal in single-cell datasets.

The procedure in Seurat is described in detail here, and improves on previous versions by directly modeling the mean-variance relationship inherent in single-cell data, and is implemented in the `FindVariableFeatures()` function. By default, Seurat returns 2,000 features per dataset. These will be used in downstream analysis, like PCA.

```
seurat_object <- FindVariableFeatures(seurat_object, selection.method = 'vst', nfeatures = 2000)
#> Finding variable features for layer counts
# Identify the 10 most highly variable genes
top_vf <- head(VariableFeatures(seurat_object), 10)
# plot variable features with and without labels
plot1 <- VariableFeaturePlot(seurat_object)
```

```
plot2 <- LabelPoints(plot1, points = top_vf, repel = TRUE)
#> When using repel, set xnudge and ynudge to 0 for optimal results
plot2
#> Warning in scale_x_log10(): log-10 transformation
#> introduced infinite values.
```



## 7.2 Scaling the data

### Why do we need to do this?

Highly expressed genes can overpower the signal of other less expressed genes with equal importance. Within the same cell the assumption is that the underlying RNA content is constant. Additionally, If variables are provided in vars.to.regress, they are individually regressed against each feature, and the resulting residuals are then scaled and centered. This step allows controlling for cell cycle and other factors that may bias your clustering.

Next, we apply a linear transformation ('scaling') that is a standard pre-processing step prior to dimensional reduction techniques like PCA. The `ScaleData()` function:

- Shifts the expression of each gene, so that the mean expression across cells is 0

- Scales the expression of each gene, so that the variance across cells is 1
  - This step gives equal weight in downstream analyses, so that highly-expressed genes do not dominate
- The results of this are stored in `seurat_object$RNA@scale.data`

```
all.genes <- rownames(seurat_object)
seurat_object <- ScaleData(seurat_object, features = all.genes)
#> Centering and scaling data matrix
```

**This step takes too long! Can I make it faster?**

Scaling is an essential step in the Seurat workflow, but only on genes that will be used as input to PCA. Therefore, the default in `ScaleData()` is only to perform scaling on the previously identified variable features (2,000 by default). To do this, omit the `features` argument in the previous function call, i.e.

```
# seurat_object <- ScaleData(seurat_object)
```

Your PCA and clustering results will be unaffected. However, Seurat heatmaps (produced as shown below with `DoHeatmap()`) require genes in the heatmap to be scaled, to make sure highly-expressed genes don't dominate the heatmap. To make sure we don't leave any genes out of the heatmap later, we are scaling all genes in this tutorial.

### How can I remove unwanted sources of variation, as in Seurat v2?

In Seurat v2 we also use the `ScaleData()` function to remove unwanted sources of variation from a single-cell dataset. For example, we could ‘regress out’ heterogeneity associated with (for example) cell cycle stage, or mitochondrial contamination. These features are still supported in `ScaleData()` in Seurat v3, i.e.:

```
# seurat_object <- ScaleData(seurat_object, vars.to.regress = 'percent.mt')
```

However, particularly for advanced users who would like to use this functionality, the Seurat authors strongly recommend the use of their new normalization workflow, `SCTransform()`. The method is described in their paper, with a separate vignette using Seurat v3 here. As with `ScaleData()`, the function `SCTransform()` also includes a `vars.to.regress` parameter.

---



# Chapter 8

# Dimensionality reduction

## Why do we need to do this?

Imagine each gene represents a dimension - or an axis on a plot. We could plot the expression of two genes with a simple scatterplot. But a genome has thousands of genes - how do you collate all the information from each of those genes in a way that allows you to visualise it in a 2 dimensional image. This is where dimensionality reduction comes in, we calculate meta-features that contains combinations of the variation of different genes. From thousands of genes, we end up with 10s of meta-features

## 8.1 Perform linear dimensional reduction

Next we perform PCA on the scaled data. By default, only the previously determined variable features are used as input, but can be defined using `features` argument if you wish to choose a different subset.

```
seurat_object <- RunPCA(seurat_object, features = VariableFeatures(object = seurat_object))
#> PC_ 1
#> Positive: CCR7, TRAT1, CREM, ALOX5AP, NKG7, TSC22D3, CST7, PASK, GPR171, CD8A
#>           CD8B, ADTRP, SVIP, PRF1, MYC, NOP58, TTC39C, SESN3, CMTM8, C14orf1
#>           GRAP2, GZMA, TARBP1, KLRD1, CD320, GCHFR, AMICA1, TUBA4A, DDIT4, NOB1
#> Negative: TYROBP, SOD2, TIMP1, TYMP, ANXA5, LGALS3, KYNU, FCN1, LYZ, APOBEC3A
#>           CD68, NPC2, S100A11, CTSL, MAFB, HLA-DRA, SDCBP, S100A10, PLAUR, GSTO1
#>           IL4I1, IDO1, PILRA, LILRB4, S100A9, MS4A7, FGL2, CXCL11, HLA-DRB1, C3AR1
#> PC_ 2
#> Positive: CD14, S100A8, PID1, CD9, GPX1, THBS1, PLAUR, C19orf59, OSM, CTB-61M7.2
#>           MGST1, S100A9, GAPDH, C5AR1, SLC7A11, ATP6VOB, PPIF, CXCL2, TGFB1, PFN1
#>           LIMS1, OLR1, PLIN2, TIMP1, COTL1, CYP1B1, PDLM7, SLC11A1, CYP27A1, CEBPB
```

```

#> Negative: IFIT3, MX1, TNFSF10, IFIT2, IFI6, RSAD2, OAS1, CXCL11, MT2A, IRF7
#>      IFITM3, OASL, GBP1, IDO1, PLSCR1, DDX58, CMPK2, APOBEC3A, FAM26F, BST2
#>      HES4, IFIH1, RABGAP1L, IL27, VAMP5, SERPING1, GMPR, SPATS2L, IRG1, IL4I1
#> PC_ 3
#> Positive: HLA-DQA1, CD83, HLA-DQB1, HLA-DRB1, HLA-DRA, HERPUD1, HSP90AB1,
#>      PKIB, TCF4, FABP5, BANK1, HSPD1, CLIC2, CD79B, FSCN1, HSPH1, CMTM6
#>      SQLE, TNFRSF13B, CD40, ALDH2, LY9, NME1, CKS2, HSP90AA1, HAPLN3, IGLL5
#> Negative: ANXA1, NKG7, PRF1, GZMA, KLRD1, MT2A, S100A8, S100A9, OASL, CD300E
#>      CST7, C3AR1, CD8A, TYROBP, CD14, S100A12, S100A6, FCGR3A, CTSL, FCN1
#>      MAFB, GCHFR, KLRC1, S100A11, IFI6, C5AR1, AQP9, FPR1, C19orf59, CD8B
#> PC_ 4
#> Positive: CCR7, ADTRP, TRAT1, MYC, CMTM8, PASK, TARBP1, CTSL, SOCS3, S100A9
#>      S100A8, EMP3, SGTB, TSC22D3, FBLN7, SESN3, GBP1, NEXN, NPC2, MPRIP
#>      HSP90AB1, IL27, S100A12, CCR1, PPA1, FCN1, SOD2, RSAD2, GPR171, HSPD1
#> Negative: NKG7, CST7, PRF1, KLRD1, GZMA, ID2, KLRC1, TNFRSF18, RAMP1, IGFBP7
#>      ALOX5AP, CD8A, GCHFR, GNG2, GAPDH, FCGR3A, XCL2, ANXA1, PRR5L, OASL
#>      RAB27A, HAVCR2, EIF4EBP1, PKIB, RHOC, GZMK, LINC00996, ADAM8, GSN, BST2
#> PC_ 5
#> Positive: FCGR3A, MS4A4A, MS4A7, CXCL16, PPM1N, SMPDL3A, AIF1, SERPINA1, ADA, CDKN
#>      CH25H, C3AR1, PLAC8, IL3RA, PILRA, CFD, CLEC12A, VMP1, FGL2, VNN2
#>      FCGR3B, MTMR11, C1QA, MAPKAPK3, LILRB2, COTL1, FAM26F, FPR2, IFNGR2, IL15
#> Negative: S100A9, MGST1, SLC7A11, S100A8, P2RY6, CCR5, LYZ, FPR3, RSAD2, FABP5
#>      SDS, EMP1, IFI6, CCR1, DHRS9, PRF1, CSTB, MX1, LILRB4, SPHK1
#>      TGFBI, ANXA1, IDO1, CXCL2, CCNA1, HSP90B1, NKG7, CMPK2, S100A12, LPXN

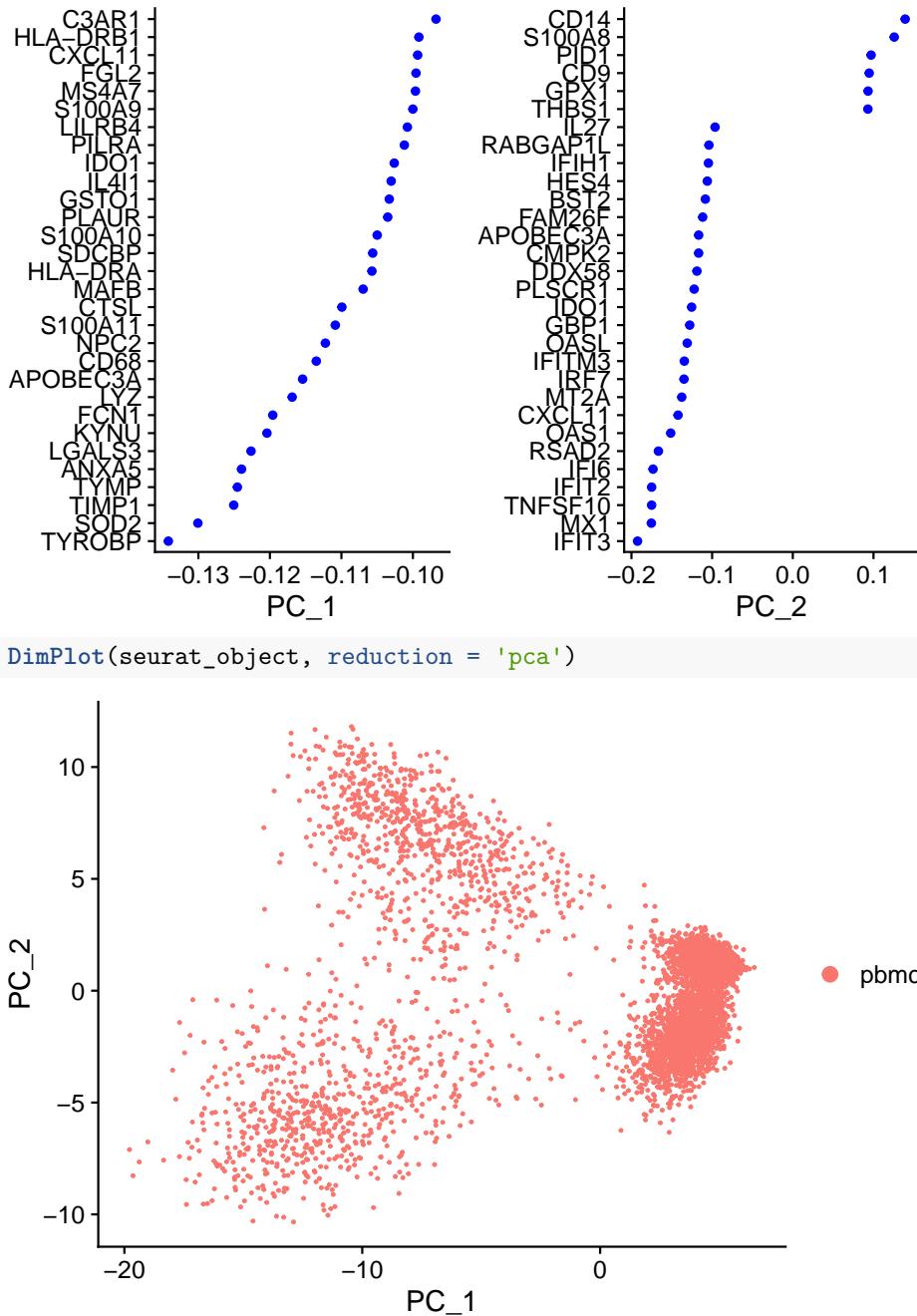
```

Seurat provides several useful ways of visualizing both cells and features that define the PCA, including `VizDimReduction()`, `DimPlot()`, and `DimHeatmap()`

```

# Examine and visualize PCA results a few different ways
print(seurat_object$pca, dims = 1:5, nfeatures = 5)
#> PC_ 1
#> Positive: CCR7, TRAT1, CREM, ALOX5AP, NKG7
#> Negative: TYROBP, SOD2, TIMP1, TYMP, ANXA5
#> PC_ 2
#> Positive: CD14, S100A8, PID1, CD9, GPX1
#> Negative: IFIT3, MX1, TNFSF10, IFIT2, IFI6
#> PC_ 3
#> Positive: HLA-DQA1, CD83, HLA-DQB1, HLA-DRB1, HLA-DRA
#> Negative: ANXA1, NKG7, PRF1, GZMA, KLRD1
#> PC_ 4
#> Positive: CCR7, ADTRP, TRAT1, MYC, CMTM8
#> Negative: NKG7, CST7, PRF1, KLRD1, GZMA
#> PC_ 5
#> Positive: FCGR3A, MS4A4A, MS4A7, CXCL16, PPM1N
#> Negative: S100A9, MGST1, SLC7A11, S100A8, P2RY6
VizDimLoadings(seurat_object, dims = 1:2, reduction = 'pca')

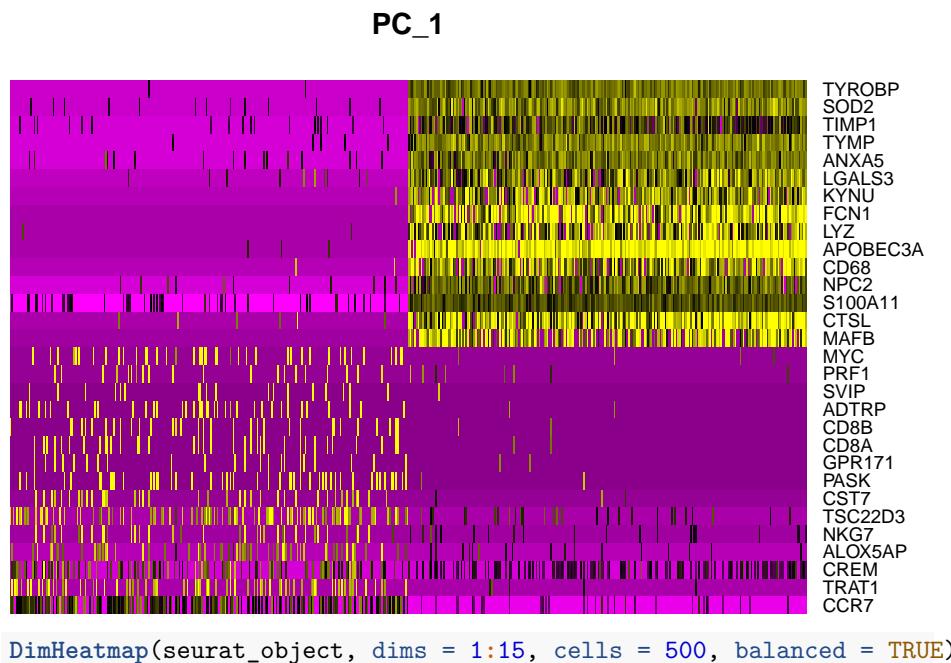
```



In particular `DimHeatmap()` allows for easy exploration of the primary sources of heterogeneity in a dataset, and can be useful when trying to decide which PCs to include for further downstream analyses. Both cells and features are ordered according to their PCA scores. Setting `cells` to a number plots the

'extreme' cells on both ends of the spectrum, which dramatically speeds plotting for large datasets. Though clearly a supervised analysis, we find this to be a valuable tool for exploring correlated feature sets.

```
DimHeatmap(seurat_object, dims = 1, cells = 500, balanced = TRUE)
```





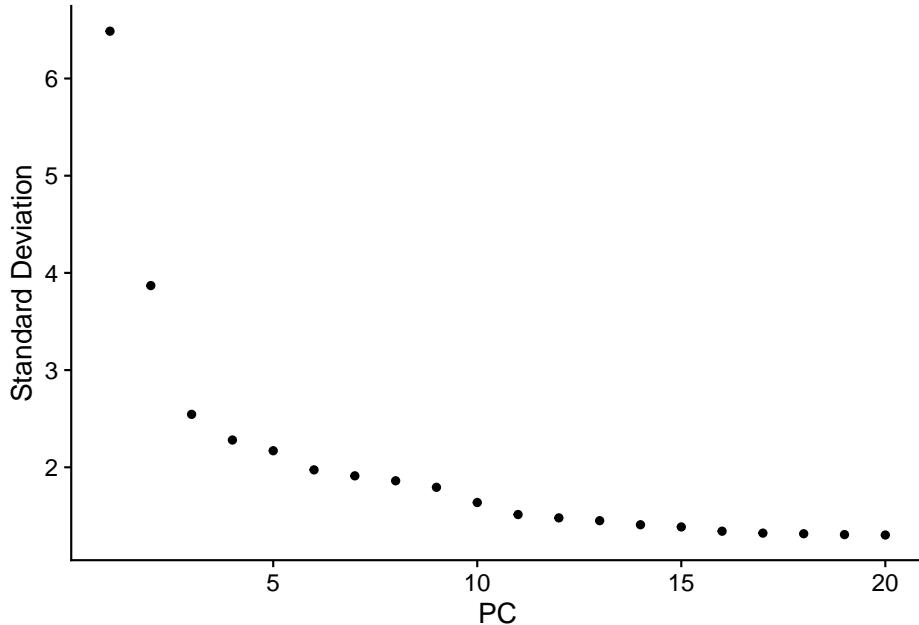
## 8.2 Determine the ‘dimensionality’ of the dataset

To overcome the extensive technical noise in any single feature for scRNA-seq data, Seurat clusters cells based on their PCA scores, with each PC essentially representing a ‘metafeature’ that combines information across a correlated feature set. The top principal components therefore represent a robust compression of the dataset. However, how many components should we choose to include? 10? 20? 100?

*Note:* The Seurat developers suggest using a JackStraw resampling test to determine this. See Macosko *et al.*, and the original `seurat_object3` vignette. We’re going to use an Elbow Plot instead here, because its much quicker.

An alternative heuristic method generates an ‘Elbow plot’: a ranking of principle components based on the percentage of variance explained by each one (`ElbowPlot()` function). In this example, we can observe an ‘elbow’ around PC9-10, suggesting that the majority of true signal is captured in the first 10 PCs.

```
ElbowPlot(seurat_object)
```



Identifying the true dimensionality of a dataset – can be challenging/uncertain for the user. We therefore suggest these three approaches to consider. The first

is more supervised, exploring PCs to determine relevant sources of heterogeneity, and could be used in conjunction with GSEA for example. The second implements a statistical test based on a random null model, but is time-consuming for large datasets, and may not return a clear PC cutoff. The third is a heuristic that is commonly used, and can be calculated instantly. In this example, all three approaches yielded similar results, but we might have been justified in choosing anything between PC 7-12 as a cutoff.

We chose 10 here, but encourage users to consider the following:

- In the original version of this vignette with the PBMC3k dataset, genes strongly associated with PCs 12 and 13 defined rare immune subsets (i.e. MZB1 is a marker for plasmacytoid DCs). However, these groups are so rare, they are difficult to distinguish from background noise for a dataset of this size without prior knowledge.
  - We encourage users to repeat downstream analyses with a different number of PCs (10, 15, or even 50!). As you will observe, the results often do not differ dramatically.
  - We advise users to err on the higher side when choosing this parameter. For example, performing downstream analyses with only 5 PCs does significantly and adversely affect results.
- 

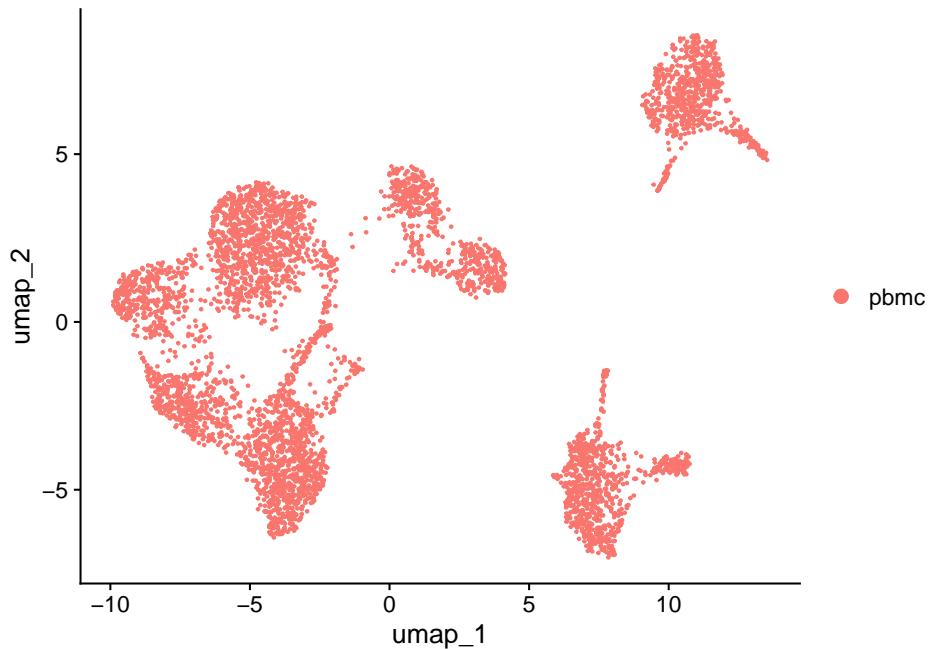
### 8.3 Run non-linear dimensional reduction (UMAP/tSNE)

Seurat offers several non-linear dimensional reduction techniques, such as tSNE and UMAP, to visualize and explore these datasets. The goal of these algorithms is to learn the underlying manifold of the data in order to place similar cells together in low-dimensional space. Cells within the graph-based clusters determined above should co-localize on these dimension reduction plots. As input to the UMAP and tSNE, we suggest using the same PCs as input to the clustering analysis.

```
seurat_object <- RunUMAP(seurat_object, dims = 1:10)
#> Warning: The default method for RunUMAP has changed from calling Python UMAP via reticulate to
#> To use Python UMAP via reticulate, set umap.method to 'umap-learn' and metric to 'correlation'
#> This message will be shown once per session
#> 23:33:13 UMAP embedding parameters a = 0.9922 b = 1.112
#> 23:33:13 Read 4877 rows and found 10 numeric columns
#> 23:33:13 Using Annoy for neighbor search, n_neighbors = 30
#> 23:33:13 Building Annoy index with metric = cosine, n_trees = 50
#> 0%   10  20  30  40  50  60  70  80  90  100%
#> [----/----/----/----/----/----/----/----/----/
#> *****/******/******/******/******/******/******/*****/
```

```
#> 23:33:13 Writing NN index file to temp file /tmp/RtmporEusV/file1ab2d9648a275d
#> 23:33:13 Searching Annoy index using 1 thread, search_k = 3000
#> 23:33:15 Annoy recall = 100%
#> 23:33:16 Commencing smooth kNN distance calibration using 1 thread with target n_ne
#> 23:33:17 Initializing from normalized Laplacian + noise (using RSpectra)
#> 23:33:17 Commencing optimization for 500 epochs, with 201470 positive edges
#> 23:33:17 Using rng type: pcg
#> 23:33:22 Optimization finished
```

```
DimPlot(seurat_object, reduction = 'umap')
```



### Challenge: PC genes

You can plot gene expression on the UMAP with the `FeaturePlot()` function.

Try out some genes that were highly weighted in the principal component analysis. How do they look?

## 8.4 Save

You can save the object at this point so that it can easily be loaded back in with `readRDS()` without having to rerun the computationally intensive steps performed above, or easily shared with collaborators.

```
saveRDS(seurat_object, file = "seurat_object_tutorial_saved.rds")
```

Tip: For faster saving and loading, try the “qs” package.



# Chapter 9

## Data set integration with Harmony

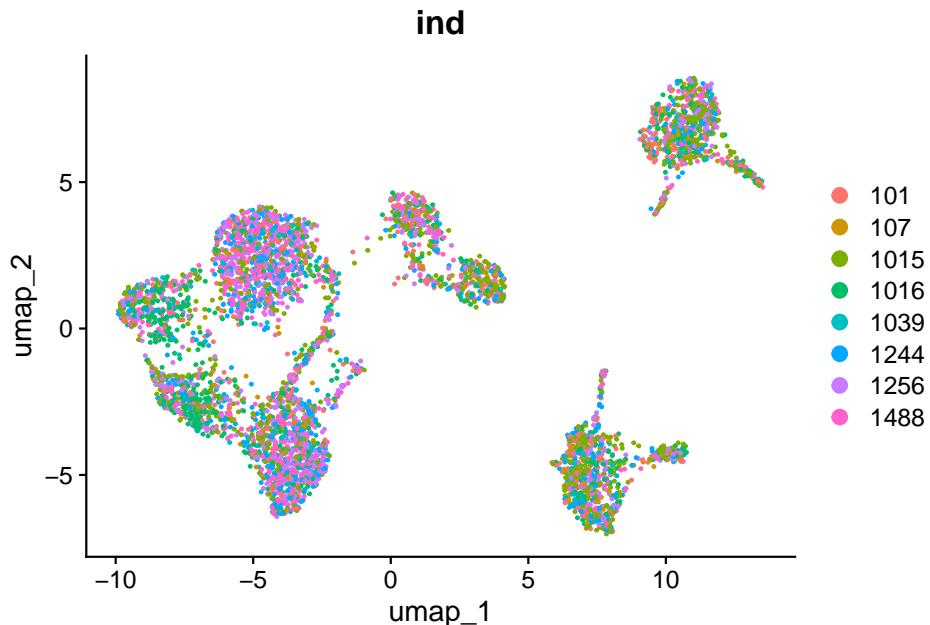
Slides

### Why do we need to do this?

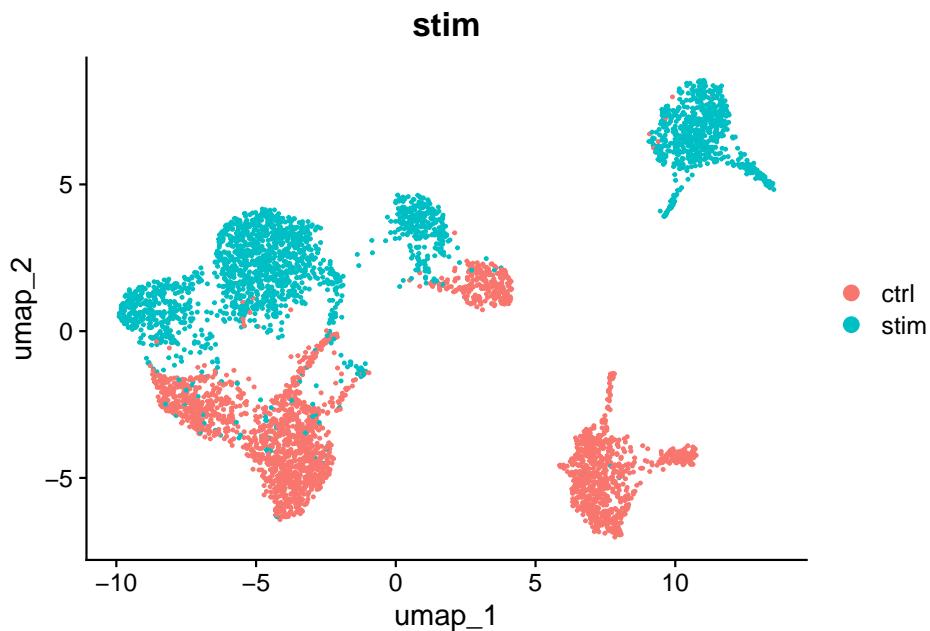
You can have data coming from different samples, batches or experiments and you will need to combine them.

- `ind` identifies a cell as coming from one of 8 individuals.
- `stim` identifies a cell as control or stimulated with IFN-beta.
- `cell` contains the cell types identified by the creators of this data set.
- `multiplets` classifies cells as singlet or doublet.

```
DimPlot(seurat_object, reduction="umap", group.by="ind")
```



```
DimPlot(seurat_object, reduction="umap", group.by="stim")
```



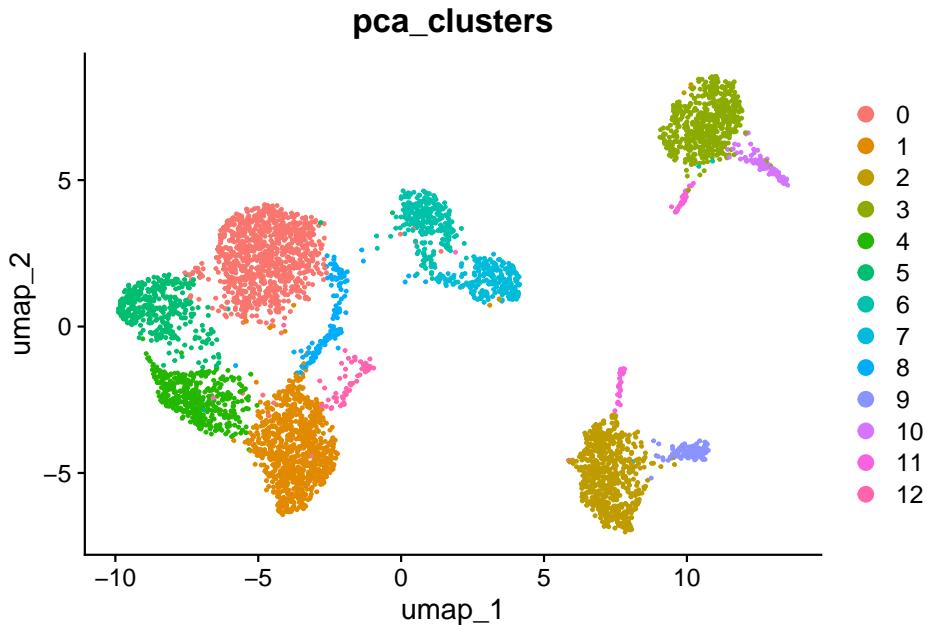
```
seurat_object<- FindNeighbors(seurat_object, reduction="pca", dims=1:10)
#> Computing nearest neighbor graph
#> Computing SNN
```

```

seurat_object <- FindClusters(seurat_object, resolution=0.5)
#> Modularity Optimizer version 1.3.0 by Ludo Waltman and Nees Jan van Eck
#>
#> Number of nodes: 4877
#> Number of edges: 171734
#>
#> Running Louvain algorithm...
#> Maximum modularity in 10 random starts: 0.9108
#> Number of communities: 13
#> Elapsed time: 0 seconds
seurat_object$pca_clusters <- seurat_object$seurat_clusters

DimPlot(seurat_object, reduction="umap", group.by="pca_clusters")

```



There is a big difference between unstimulated and stimulated cells. This has split cells of the same type into pairs of clusters. If the difference was simply uniform, we could regress it out (e.g. using `ScaleData(..., vars.to.regress="stim")`). However, as can be seen in the PCA plot, the difference is not uniform and we need to do something cleverer.

We will use Harmony, which can remove non-uniform effects. We will try to remove both the small differences between individuals and the large difference between the unstimulated and stimulated cells.

Harmony operates only on the PCA scores. The original gene expression levels remain unaltered.

```

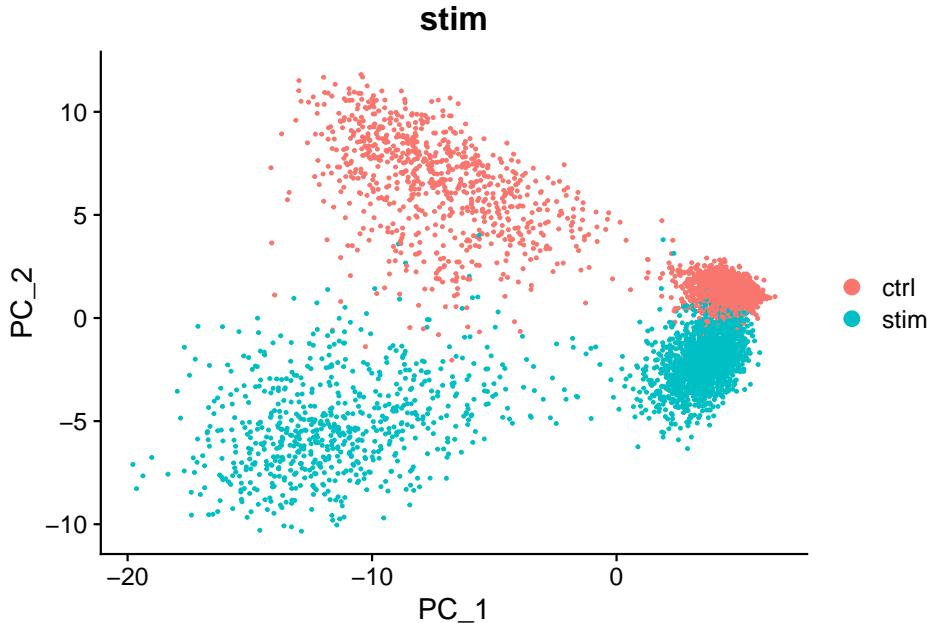
library(harmony)
#> Loading required package: Rcpp

seurat_object <- RunHarmony(seurat_object, c("stim", "ind"), reduction="pca", reduction
#> Transposing data matrix
#> Initializing state using k-means centroids initialization
#> Harmony 1/10
#> Harmony 2/10
#> Harmony 3/10
#> Harmony 4/10
#> Harmony converged after 4 iterations
#> Warning: The `slot` argument of `GetAssayData()` is deprecated as of
#> SeuratObject 5.0.0.
#> i Please use the `layer` argument instead.
#> i The deprecated feature was likely used in the Seurat
#> package.
#> Please report the issue at
#> <https://github.com/satijalab/seurat/issues>.
#> This warning is displayed once every 8 hours.
#> Call `lifecycle::last_lifecycle_warnings()` to see where
#> this warning was generated.

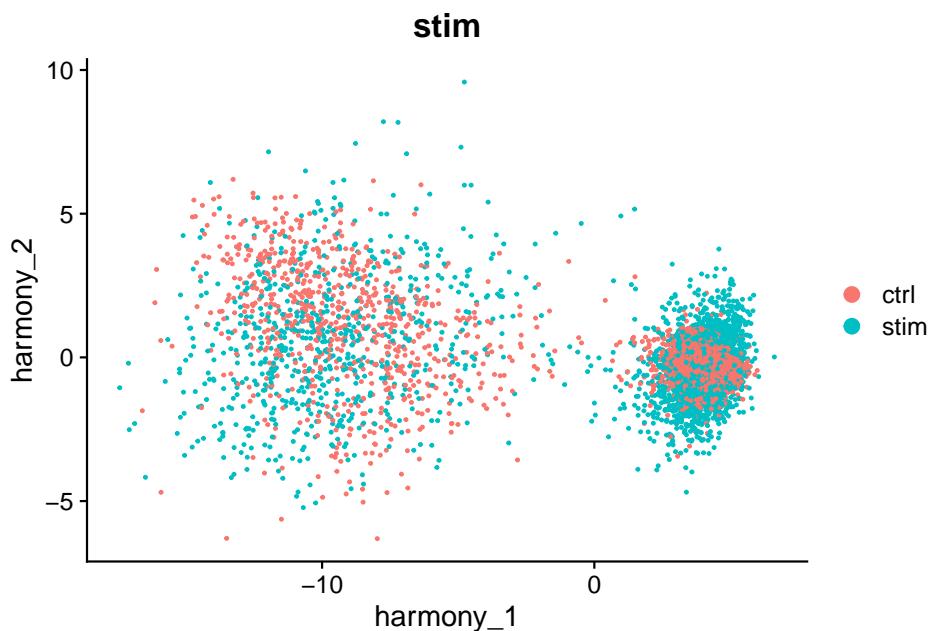
```

This has added a new set of reduced dimensions to the Seurat object, `seurat_object$harmony` which is a modified version of the existing `seurat_object$pca` reduced dimensions. The PCA plot shows a large difference between ‘ctrl’ and ‘stim’, but this has been removed in the harmony reduction.

```
DimPlot(seurat_object, reduction="pca", group.by="stim")
```



```
DimPlot(seurat_object, reduction="harmony", group.by="stim")
```



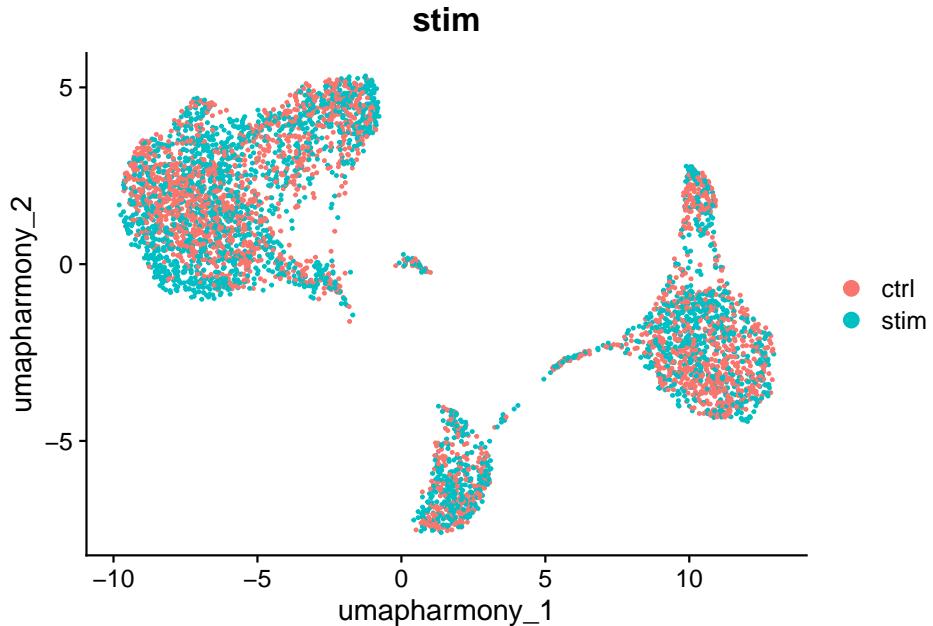
We can use `harmony` the same way we used the `pca` reduction to compute a UMAP layout or to find clusters.

```

seurat_object <- RunUMAP(seurat_object, reduction="harmony", dims=1:10, reduction.name=
#> 23:33:49 UMAP embedding parameters a = 0.9922 b = 1.112
#> 23:33:49 Read 48771 rows and found 10 numeric columns
#> 23:33:49 Using Annoy for neighbor search, n_neighbors = 30
#> 23:33:49 Building Annoy index with metric = cosine, n_trees = 50
#> 0%   10  20  30  40  50  60  70  80  90  100%
#> [----|----|----|----|----|----|----|----|----|----|
#> ****|*****|*****|*****|*****|*****|*****|*****|*****|*****|
#> 23:33:49 Writing NN index file to temp file /tmp/RtmporEusV/file1ab2d9146582e1
#> 23:33:49 Searching Annoy index using 1 thread, search_k = 3000
#> 23:33:51 Annoy recall = 100%
#> 23:33:52 Commencing smooth kNN distance calibration using 1 thread with target n_ne
#> 23:33:53 Initializing from normalized Laplacian + noise (using RSpectra)
#> 23:33:53 Commencing optimization for 500 epochs, with 203710 positive edges
#> 23:33:53 Using rng type: pcg
#> 23:33:58 Optimization finished

```

```
DimPlot(seurat_object, reduction="umap_harmony", group.by="stim")
```



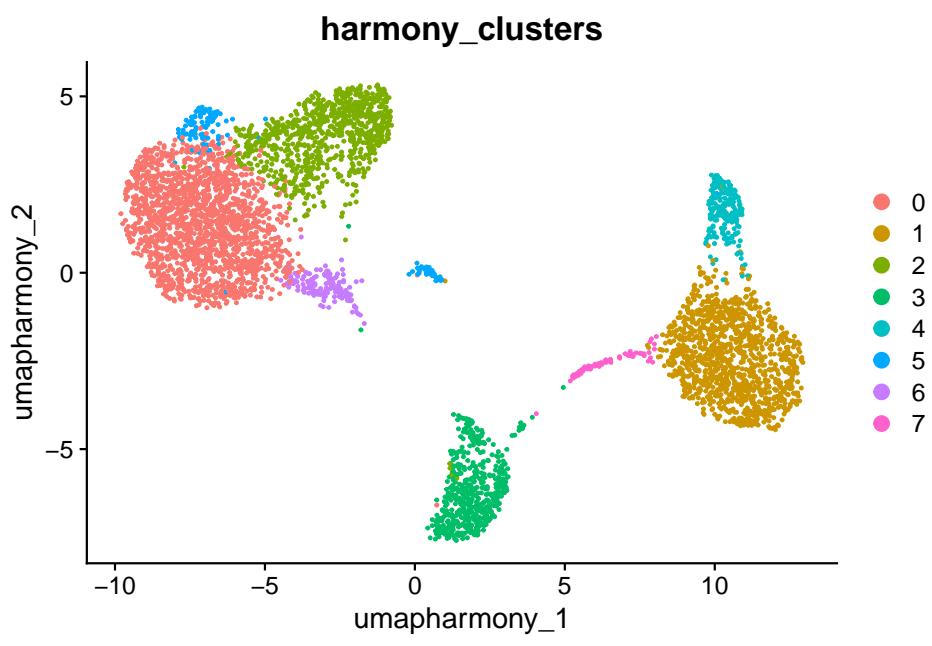
```

seurat_object <- FindNeighbors(seurat_object, reduction="harmony", dims=1:10)
#> Computing nearest neighbor graph
#> Computing SNN
seurat_object <- FindClusters(seurat_object, resolution=0.5)
#> Modularity Optimizer version 1.3.0 by Ludo Waltman and Nees Jan van Eck
#>

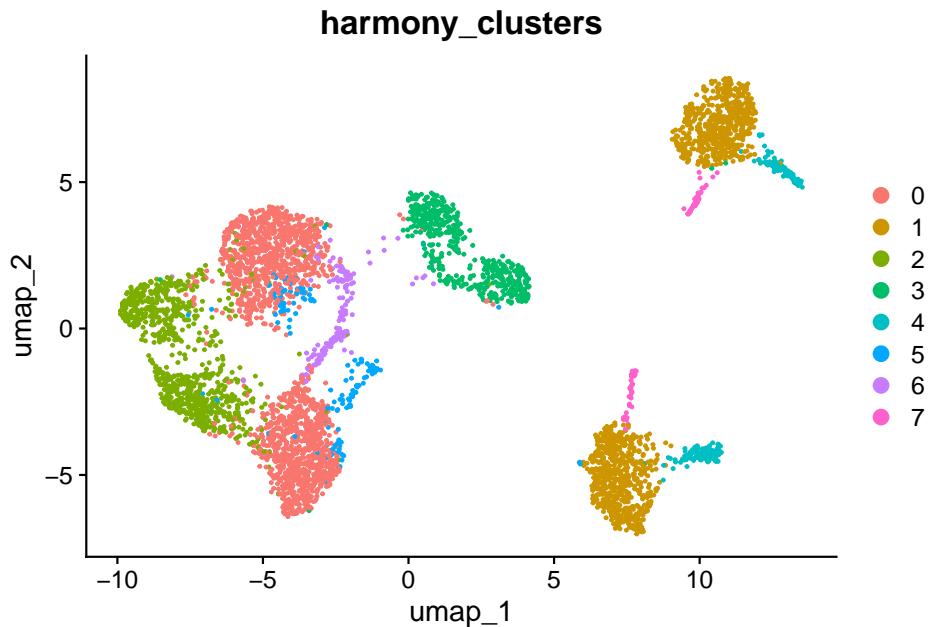
```

```
#> Number of nodes: 4877
#> Number of edges: 169427
#>
#> Running Louvain algorithm...
#> Maximum modularity in 10 random starts: 0.8660
#> Number of communities: 8
#> Elapsed time: 0 seconds
seurat_object$harmony_clusters <- seurat_object$seurat_clusters

DimPlot(seurat_object, reduction="umap_harmony", group.by="harmony_clusters")
```



```
DimPlot(seurat_object, reduction="umap", group.by="harmony_clusters")
```



Having found a good set of clusters, we would usually perform differential expression analysis on the original data and include batches/runs/individuals as predictors in the linear model. In this example we could now compare unstimulated and stimulated cells within each cluster. A particularly nice statistical approach that is possible here would be to convert the counts to pseudo-bulk data for the eight individuals, and then apply a bulk RNA-Seq differential expression analysis method. However there is still the problem that unstimulated and stimulated cells were processed in separate batches.

# Chapter 10

## Clustering

Slides

### Why do we need to do this?

Clustering the cells will allow you to visualise the variability of your data, can help to segregate cells into cell types.

### 10.1 Cluster cells

Seurat v3 applies a graph-based clustering approach, building upon initial strategies in (Macosko *et al.*). Importantly, the *distance metric* which drives the clustering analysis (based on previously identified PCs) remains the same. However, our approach to partitioning the cellular distance matrix into clusters has dramatically improved. Our approach was heavily inspired by recent manuscripts which applied graph-based clustering approaches to scRNA-seq data [SNN-Cliq, Xu and Su, Bioinformatics, 2015] and CyTOF data [PhenoGraph, Levine *et al.*, Cell, 2015]. Briefly, these methods embed cells in a graph structure - for example a K-nearest neighbor (KNN) graph, with edges drawn between cells with similar feature expression patterns, and then attempt to partition this graph into highly interconnected ‘quasi-cliques’ or ‘communities’.

As in PhenoGraph, we first construct a KNN graph based on the euclidean distance in PCA space, and refine the edge weights between any two cells based on the shared overlap in their local neighborhoods (Jaccard similarity). This step is performed using the `FindNeighbors()` function, and takes as input the previously defined dimensionality of the dataset (first 10 PCs).

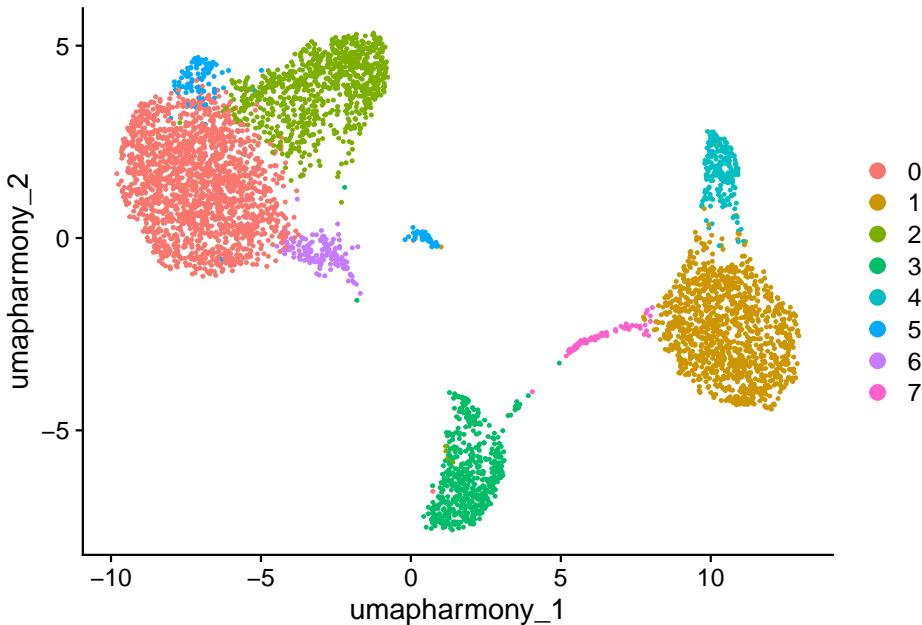
To cluster the cells, we next apply modularity optimization techniques such as the Louvain algorithm (default) or SLM [SLM, Blondel *et al.*, Journal of Statist-

tical Mechanics], to iteratively group cells together, with the goal of optimizing the standard modularity function. The `FindClusters()` function implements this procedure, and contains a resolution parameter that sets the ‘granularity’ of the downstream clustering, with increased values leading to a greater number of clusters. We find that setting this parameter between 0.4-1.2 typically returns good results for single-cell datasets of around 3K cells. Optimal resolution often increases for larger datasets. The clusters can be found using the `Idents()` function.

```
seurat_object <- FindNeighbors(seurat_object, dims = 1:10, reduction = "harmony")
#> Computing nearest neighbor graph
#> Computing SNN
seurat_object <- FindClusters(seurat_object, resolution = 0.5)
#> Modularity Optimizer version 1.3.0 by Ludo Waltman and Nees Jan van Eck
#>
#> Number of nodes: 4877
#> Number of edges: 169427
#>
#> Running Louvain algorithm...
#> Maximum modularity in 10 random starts: 0.8660
#> Number of communities: 8
#> Elapsed time: 0 seconds
# Look at cluster IDs of the first 5 cells
head(Idents(seurat_object), 5)
#> AGGGCGCTATTCC-1 GGAGACGATTGTT-1 CACCGTTGTCGTAG-1
#>           1           5           4
#> TATCGTACACGCAT-1 TACGAGACCTATT-1
#>           3           0
#> Levels: 0 1 2 3 4 5 6 7
```

Check out the clusters.

```
DimPlot(seurat_object, reduction = "umap_harmony")
```



```
# Equivalent to
# DimPlot(seurat_object, reduction="umap", group.by="seurat_clusters")
# DimPlot(seurat_object, reduction="umap", group.by="RNA_snn_res.0.5")
```

---

### Challenge: Try different cluster settings

Run `FindNeighbours` and `FindClusters` again, with a different number of dimensions or with a different resolution. Examine the resulting clusters using `DimPlot`.

To maintain the flow of this tutorial, please put the output of this exploration in a different variable, such as `seurat_object2`!

## 10.2 Choosing a cluster resolution

Its a good idea to try different resolutions when clustering to identify the variability of your data.

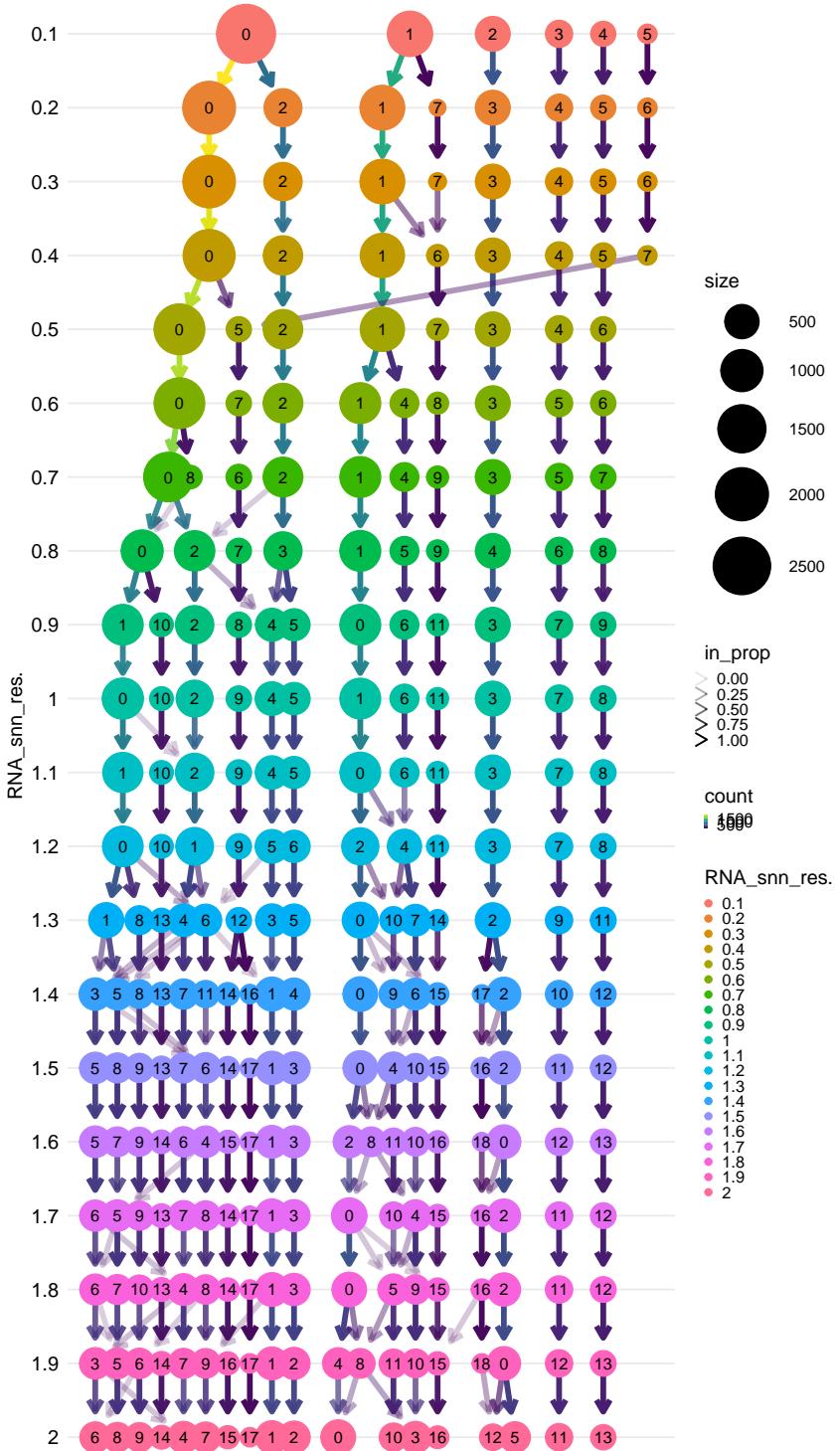
```
resolution = 2
seurat_object <- FindClusters(seurat_object, resolution = seq(0.1, resolution, 0.1))

# the different clustering created
names(seurat_object@meta.data)
```

```
# Look at cluster IDs of the first 5 cells
head(Idents(seurat_object), 5)
```

Plot a clustree to decide how many clusters you have and what resolution capture them.

```
library(clustree)
#> Loading required package: ggraph
#>
#> Attaching package: 'ggraph'
#> The following object is masked from 'package:sp':
#>
#>     geometry
clustree(seurat_object, prefix = "RNA_snn_res.", show_axis=TRUE) + theme(legend.key.size
```

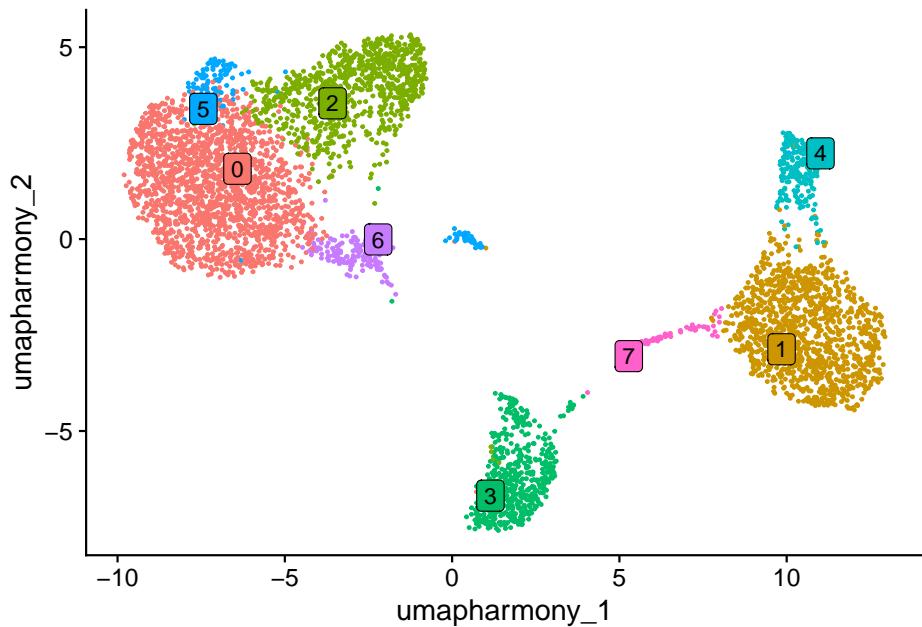


Name cells with the corresponding cluster name at the resolution you pick. This case we are happy with 0.5.

```
# The name of the cluster is prefixed with 'RNA_snn_res' and the number of the resolution
Idents(seurat_object) <- seurat_object$RNA_snn_res.0.5
```

Plot the UMAP with colored clusters with Dimplot

```
DimPlot(seurat_object, reduction = "umap_harmony", label = TRUE, repel = TRUE, label.b
```



# Chapter 11

## Cluster Markers

### Why do we need to do this?

Single cell data helps to segregate cell types. Use markers to identify cell types.  
warning: In this example the cell types/markers are well known and making this step easy, but in reality this step needs the experts curation.

### 11.1 Finding differentially expressed features (cluster biomarkers)

Seurat can help you find markers that define clusters via differential expression. By default, it identifies positive and negative markers of a single cluster (specified in `ident.1`), compared to all other cells. `FindAllMarkers()` automates this process for all clusters, but you can also test groups of clusters vs. each other, or against all cells.

The `min.pct` argument requires a feature to be detected at a minimum percentage in either of the two groups of cells, and the `thresh.test` argument requires a feature to be differentially expressed (on average) by some amount between the two groups. You can set both of these to 0, but with a dramatic increase in time - since this will test a large number of features that are unlikely to be highly discriminatory. As another option to speed up these computations, `max.cells.per.ident` can be set. This will downsample each identity class to have no more cells than whatever this is set to. While there is generally going to be a loss in power, the speed increases can be significant and the most highly differentially expressed features will likely still rise to the top.

```
# find all markers of cluster 2
cluster2.markers <- FindMarkers(seurat_object, ident.1 = 2, min.pct = 0.25)
```

```

head(cluster2.markers, n = 5)
#>      p_val avg_log2FC pct.1 pct.2 p_val_adj
#> NKG7    0   5.760837 0.802 0.055      0
#> CCL5    0   4.490605 0.749 0.106      0
#> GNLY    0   6.161437 0.663 0.048      0
#> GZMB    0   5.257738 0.652 0.037      0
#> CST7    0   4.782589 0.624 0.045      0
# find all markers distinguishing cluster 5 from clusters 0 and 3
cluster5.markers <- FindMarkers(seurat_object, ident.1 = 5, ident.2 = c(0, 3), min.pct = 0.25)
head(cluster5.markers, n = 5)
#>      p_val avg_log2FC pct.1 pct.2 p_val_adj
#> PF4   9.652092e-191   7.934396 0.435 0.005 3.439523e-186
#> GNG11 5.518953e-171   7.336953 0.429 0.008 1.966679e-166
#> CLU   4.848919e-126   8.077343 0.292 0.003 1.727912e-121
#> PPBP  6.404717e-108   6.987740 0.442 0.030 2.282321e-103
#> NCOA4 2.574636e-52   4.363135 0.351 0.046 9.174717e-48
# find markers for every cluster compared to all remaining cells, report only the positive ones
seurat_object.markers <- FindAllMarkers(seurat_object, only.pos = TRUE, min.pct = 0.25)
#> Calculating cluster 0
#> Calculating cluster 1
#> Calculating cluster 2
#> Calculating cluster 3
#> Calculating cluster 4
#> Calculating cluster 5
#> Calculating cluster 6
#> Calculating cluster 7
seurat_object.markers %>% group_by(cluster) %>% slice_max(n = 2, order_by = avg_log2FC)
#> # A tibble: 16 x 7
#> # Groups:   cluster [8]
#>      p_val avg_log2FC pct.1 pct.2 p_val_adj cluster gene
#>      <dbl>     <dbl> <dbl> <dbl>     <dbl> <fct>  <chr>
#> 1  5.85e-192    2.07  0.611  0.224  2.09e-187 0     SELL
#> 2  1.04e-224    2.04  0.654  0.22   3.71e-220 0     LTB
#> 3  7.51e-269    7.10  0.354  0.012  2.68e-264 1     CCL7
#> 4  0             7.03  0.673  0.02   0         1     S100A8
#> 5  0             6.16  0.663  0.048  0         2     GNLY
#> 6  3.91e-259    6.12  0.333  0.011  1.39e-254 2     GZMA
#> 7  0             6.26  0.444  0.011  0         3     MS4A1
#> 8  0             6.07  0.658  0.019  0         3     CD79A
#> 9  0             6.71  0.833  0.029  0         4     VMO1
#> 10 0            5.19  0.742  0.036  0         4     MS4A4A
#> 11 7.61e-253    7.39  0.435  0.011  2.71e-248 5     PF4
#> 12 9.17e-110    6.53  0.442  0.042  3.27e-105 5     PPBP
#> 13 5.78e- 75    4.65  0.253  0.018  2.06e- 70 6     NR4A2
#> 14 2.99e-119    4.57  0.604  0.077  1.07e-114 6     HSPH1

```

### 11.1. FINDING DIFFERENTIALLY EXPRESSED FEATURES (CLUSTER BIOMARKERS)59

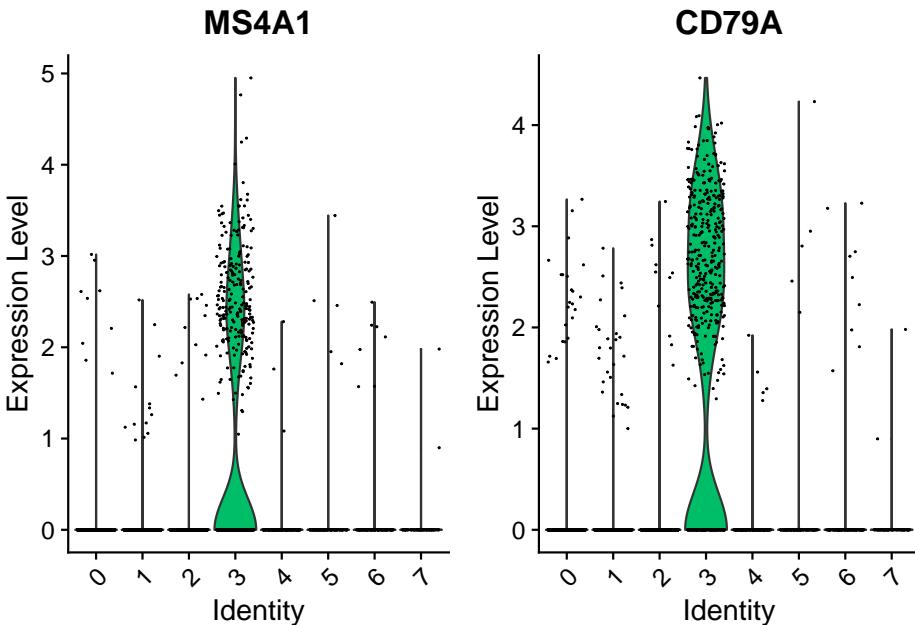
```
#> 15 0          8.82 0.448 0.001 0      7      CCL22
#> 16 1.61e-224 8.23 0.253 0.001 5.74e-220 7      CALCR
```

Seurat has several tests for differential expression which can be set with the `test.use` parameter (see our DE vignette for details). For example, the ROC test returns the ‘classification power’  $\text{abs}(\text{AUC}-0.5)*2$  for any individual marker, ranging from 0 = random to 1 = perfect.

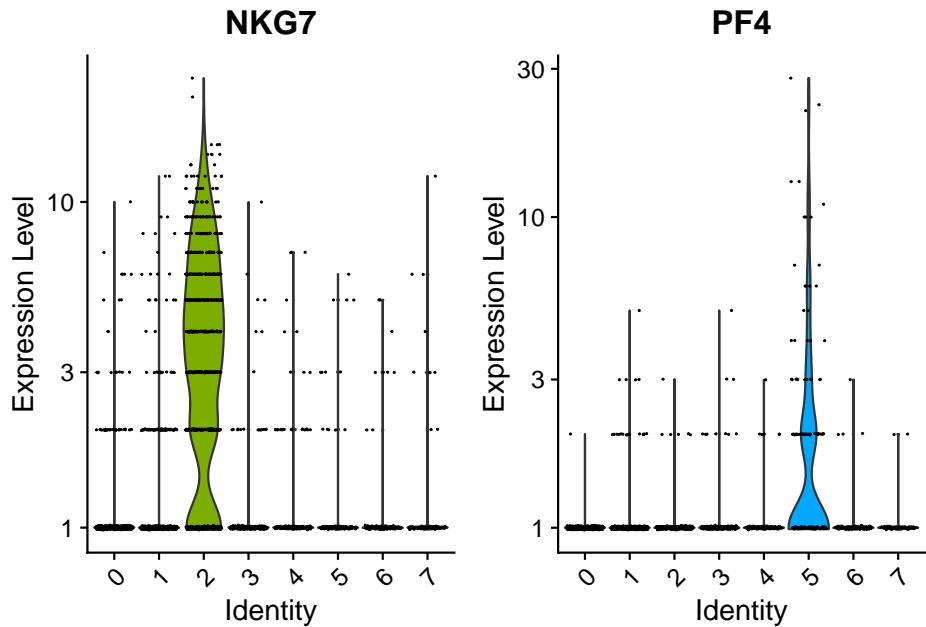
```
cluster0.markers <- FindMarkers(seurat_object, ident.1 = 0, logfc.threshold = 0.25, test.use = "r")
```

We include several tools for visualizing marker expression. `VlnPlot()` (shows expression probability distributions across clusters), and `FeaturePlot()` (visualizes feature expression on a tSNE or PCA plot) are our most commonly used visualizations. We also suggest exploring `RidgePlot()`, `CellScatter()`, and `DotPlot()` as additional methods to view your dataset.

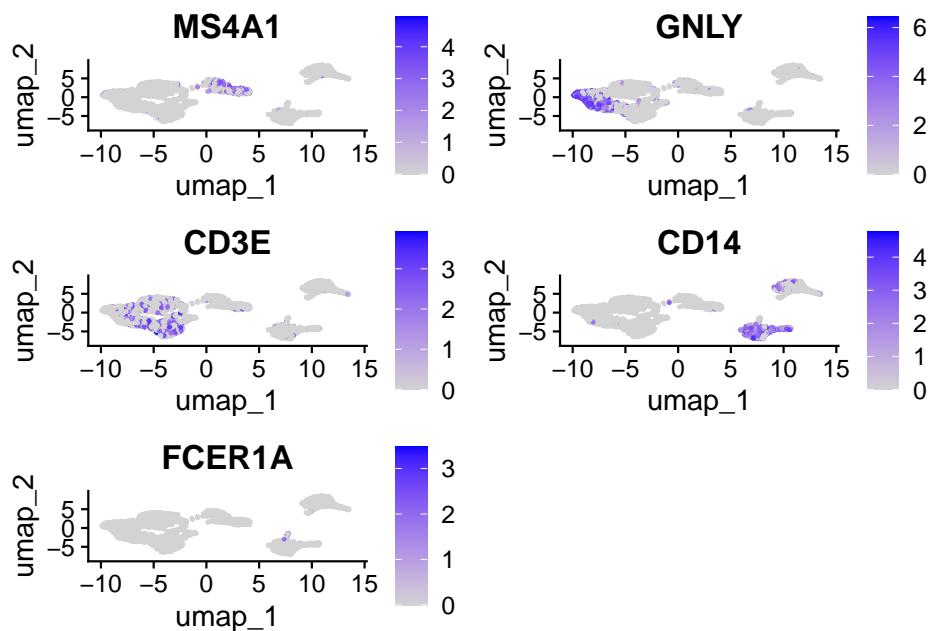
```
VlnPlot(seurat_object, features = c("MS4A1", "CD79A"))
```



```
# you can plot raw counts as well
VlnPlot(seurat_object, features = c("NKG7", "PF4"), slot = 'counts', log = TRUE)
#> Warning: The `slot` argument of `VlnPlot()` is deprecated as of
#> Seurat 5.0.0.
#> i Please use the `layer` argument instead.
#> This warning is displayed once every 8 hours.
#> Call `lifecycle::last_lifecycle_warnings()` to see where
#> this warning was generated.
```

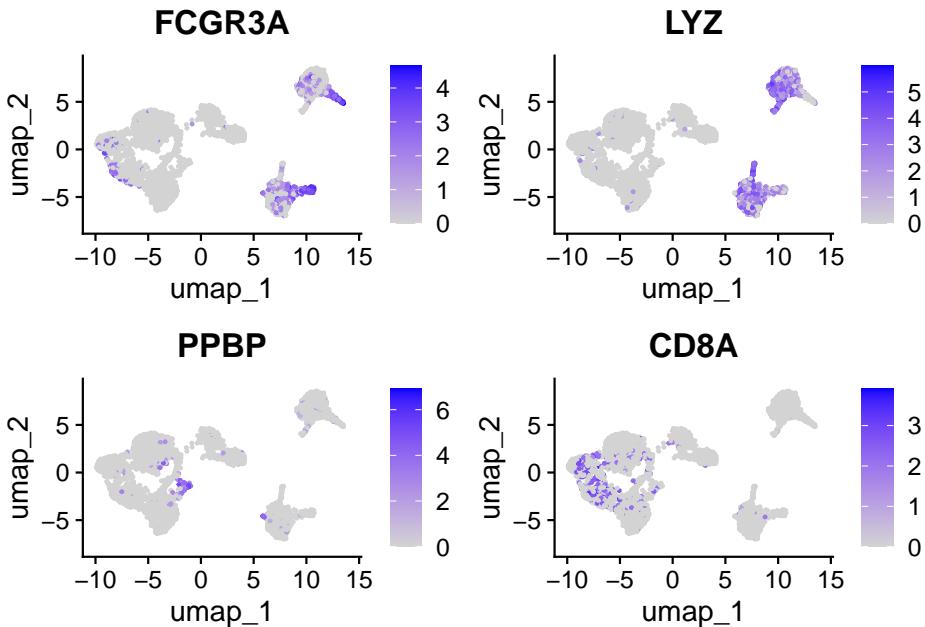


```
FeaturePlot(seurat_object, features = c("MS4A1", "GNLY", "CD3E", "CD14", "FCER1A"))
```



```
FeaturePlot(seurat_object, features = c("FCGR3A", "LYZ", "PPBP", "CD8A"))
```

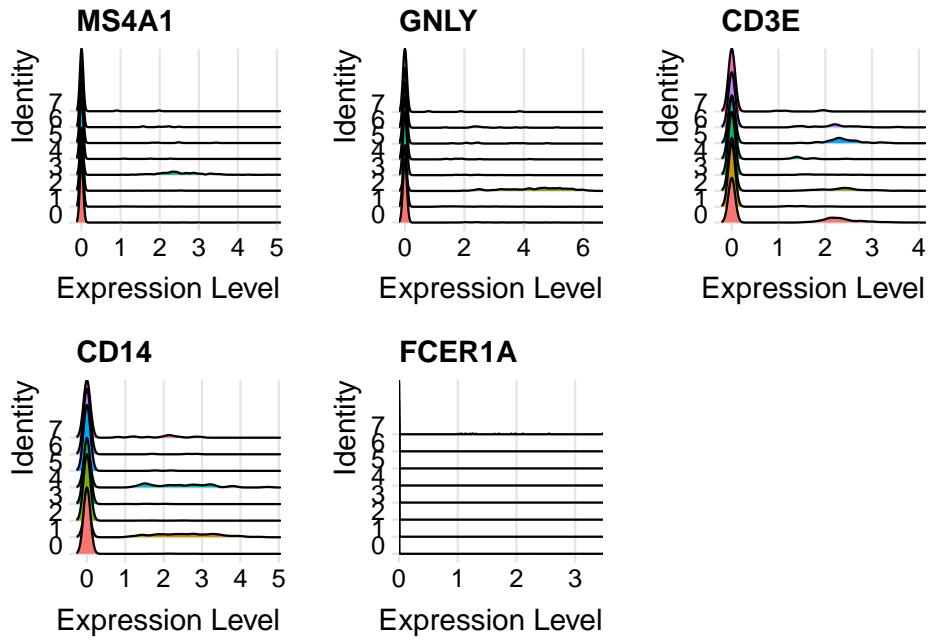
### 11.1. FINDING DIFFERENTIALLY EXPRESSED FEATURES (CLUSTER BIOMARKERS)61



### Other useful plots

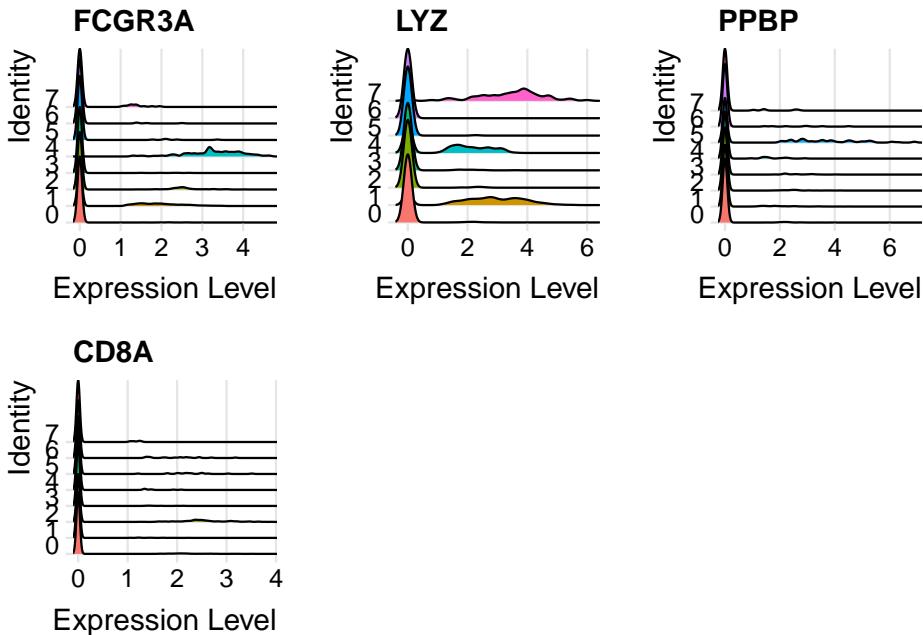
These are ridgeplots, cell scatter plots and dotplots. Replace FeaturePlot with the other functions.

```
RidgePlot(seurat_object, features = c("MS4A1", "GNLY", "CD3E", "CD14", "FCER1A"))
#> Picking joint bandwidth of 0.0427
#> Picking joint bandwidth of 0.0656
#> Picking joint bandwidth of 0.0727
#> Picking joint bandwidth of 0.0901
#> Picking joint bandwidth of 3e-06
```



```
RidgePlot(seurat_object, features = c("FCGR3A", "LYZ", "PPBP", "CD8A"))
#> Picking joint bandwidth of 0.0547
#> Picking joint bandwidth of 0.136
#> Picking joint bandwidth of 0.0835
#> Picking joint bandwidth of 0.0342
```

### 11.1. FINDING DIFFERENTIALLY EXPRESSED FEATURES (CLUSTER BIOMARKERS)63



For CellScatter plots, will need the cell id of the cells you want to look at. You can get this from the cell metadata (`seurat_object@meta.data`).

```
head( seurat_object@meta.data )
#>           orig.ident nCount_RNA nFeature_RNA ind
#> AGGGCGCTATTCC-1 SeuratProject      2053      532 1256
#> GGAGACGATTCGTT-1 SeuratProject      881       392 1256
#> CACCGTTGTCGTAG-1 SeuratProject     3130      1005 1016
#> TATCGTACACGCAT-1 SeuratProject     1042       549 1488
#> TACGAGACCTATTTC-1 SeuratProject    2425      777 1244
#> GTACTACTCATACG-1 SeuratProject    3951      1064 1256
#>           stim      cell multiplets
#> AGGGCGCTATTCC-1 stim  CD14+ Monocytes   singlet
#> GGAGACGATTCGTT-1 stim  CD4 T cells     singlet
#> CACCGTTGTCGTAG-1 ctrl FCGR3A+ Monocytes singlet
#> TATCGTACACGCAT-1 stim      B cells     singlet
#> TACGAGACCTATTTC-1 stim  CD4 T cells     singlet
#> GTACTACTCATACG-1 ctrl FCGR3A+ Monocytes singlet
#>           percent.mt RNA_snn_res.0.5 seurat_clusters
#> AGGGCGCTATTCC-1 1.6336634          1          0
#> GGAGACGATTCGTT-1 4.8809524          5         15
#> CACCGTTGTCGTAG-1 1.0655473          4         11
#> TATCGTACACGCAT-1 3.0662710          3          5
#> TACGAGACCTATTTC-1 1.0837849          0          9
#> GTACTACTCATACG-1 0.7137395          4         11
#>           pca_clusters harmony_clusters
```

```

#> AGGGCGCTATTC-1      3       1
#> GGAGACGATTGTT-1    0       5
#> CACCGTTGTCGTAG-1   9       4
#> TATCGTACACGCAT-1   6       3
#> TACGAGACCTATTG-1   0       0
#> GTACTACTCATACG-1   9       4
#> RNA_snn_res.0.1 RNA_snn_res.0.2
#> AGGGCGCTATTC-1      1       1
#> GGAGACGATTGTT-1    0       0
#> CACCGTTGTCGTAG-1   3       4
#> TATCGTACACGCAT-1   2       3
#> TACGAGACCTATTG-1   0       0
#> GTACTACTCATACG-1   3       4
#> RNA_snn_res.0.3 RNA_snn_res.0.4
#> AGGGCGCTATTC-1      1       1
#> GGAGACGATTGTT-1    0       0
#> CACCGTTGTCGTAG-1   4       4
#> TATCGTACACGCAT-1   3       3
#> TACGAGACCTATTG-1   0       0
#> GTACTACTCATACG-1   4       4
#> RNA_snn_res.0.6 RNA_snn_res.0.7
#> AGGGCGCTATTC-1      1       1
#> GGAGACGATTGTT-1    7       6
#> CACCGTTGTCGTAG-1   5       5
#> TATCGTACACGCAT-1   3       3
#> TACGAGACCTATTG-1   0       0
#> GTACTACTCATACG-1   5       5
#> RNA_snn_res.0.8 RNA_snn_res.0.9
#> AGGGCGCTATTC-1      1       0
#> GGAGACGATTGTT-1    7       8
#> CACCGTTGTCGTAG-1   6       7
#> TATCGTACACGCAT-1   4       3
#> TACGAGACCTATTG-1   0       1
#> GTACTACTCATACG-1   6       7
#> RNA_snn_res.1 RNA_snn_res.1.1
#> AGGGCGCTATTC-1      1       0
#> GGAGACGATTGTT-1    9       9
#> CACCGTTGTCGTAG-1   7       7
#> TATCGTACACGCAT-1   3       3
#> TACGAGACCTATTG-1   0       1
#> GTACTACTCATACG-1   7       7
#> RNA_snn_res.1.2 RNA_snn_res.1.3
#> AGGGCGCTATTC-1      2       0
#> GGAGACGATTGTT-1    9       12
#> CACCGTTGTCGTAG-1   7       9

```

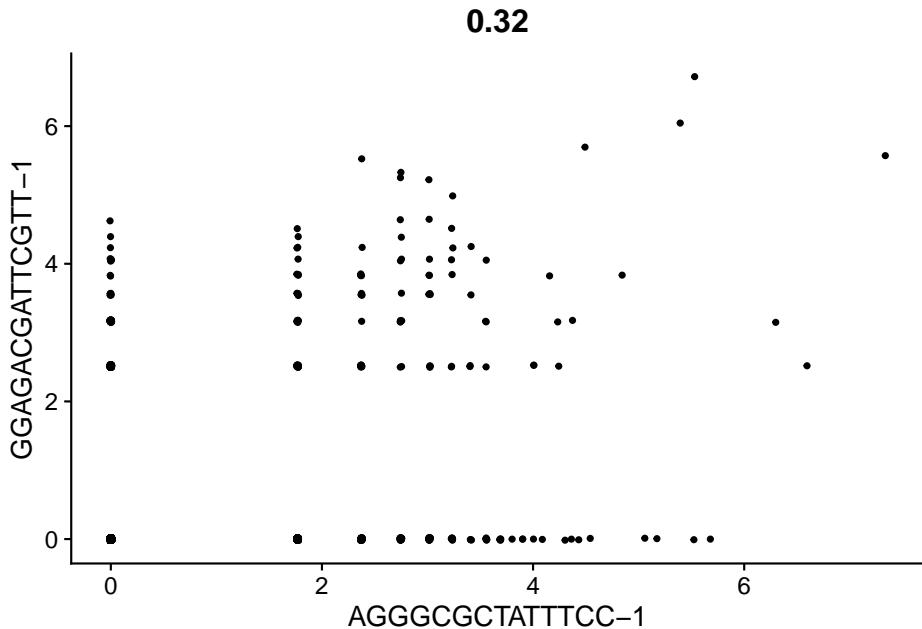
### 11.1. FINDING DIFFERENTIALLY EXPRESSED FEATURES (CLUSTER BIOMARKERS)65

```

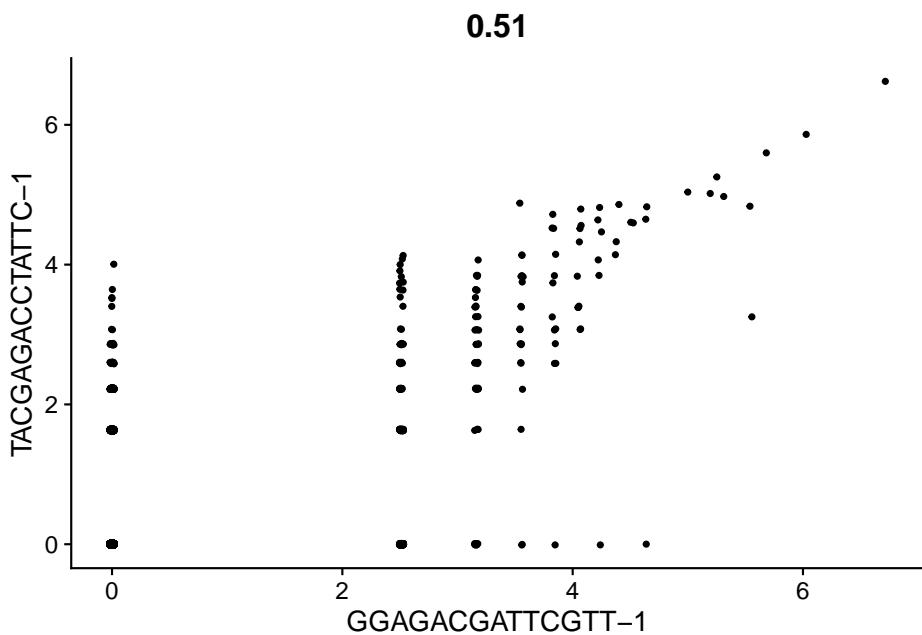
#> TATCGTACACGCAT-1          3          2
#> TACGAGACCTATTTC-1        0          8
#> GTACTACTCATACG-1         7          9
#> RNA_snn_res.1.4 RNA_snn_res.1.5
#> AGGGCGCTATTTC-1          0          0
#> GGAGACGATTGTT-1          14         14
#> CACCGTTGTCGTAG-1         10         11
#> TATCGTACACGCAT-1         2          2
#> TACGAGACCTATTTC-1        8          9
#> GTACTACTCATACG-1         10         11
#> RNA_snn_res.1.6 RNA_snn_res.1.7
#> AGGGCGCTATTTC-1          2          0
#> GGAGACGATTGTT-1          15         14
#> CACCGTTGTCGTAG-1         12         11
#> TATCGTACACGCAT-1         0          2
#> TACGAGACCTATTTC-1        9          9
#> GTACTACTCATACG-1         12         11
#> RNA_snn_res.1.8 RNA_snn_res.1.9
#> AGGGCGCTATTTC-1          0          4
#> GGAGACGATTGTT-1          14         16
#> CACCGTTGTCGTAG-1         11         12
#> TATCGTACACGCAT-1         2          0
#> TACGAGACCTATTTC-1        10         6
#> GTACTACTCATACG-1         11         12
#> RNA_snn_res.2
#> AGGGCGCTATTTC-1          0
#> GGAGACGATTGTT-1          15
#> CACCGTTGTCGTAG-1         11
#> TATCGTACACGCAT-1         5
#> TACGAGACCTATTTC-1        9
#> GTACTACTCATACG-1         11

CellScatter(seurat_object, cell1 = "AGGGCGCTATTTC-1", cell2 = "GGAGACGATTGTT-1")

```



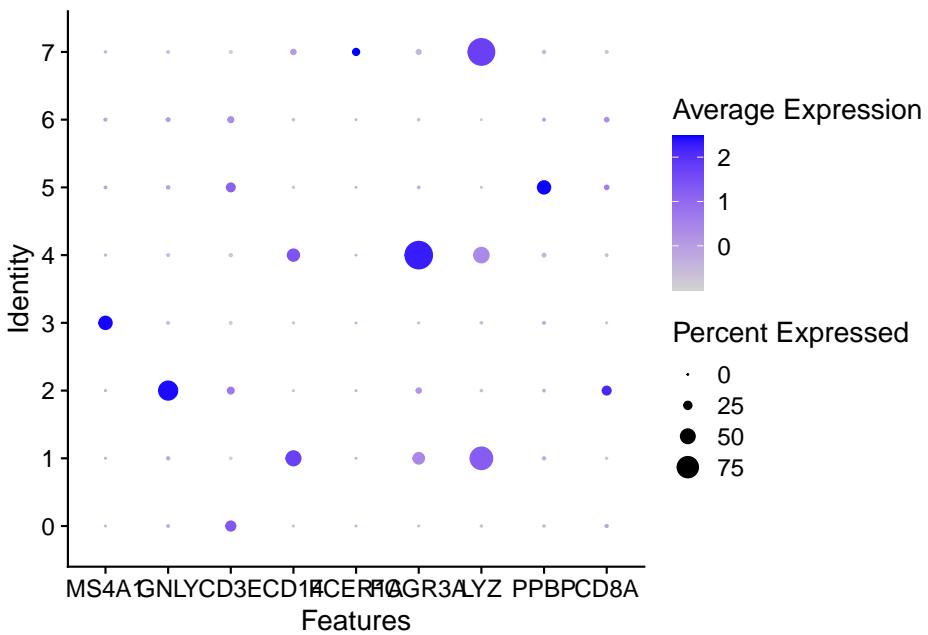
```
CellScatter(seurat_object, cell1 = "GGAGACGATTGTT-1", cell2 = "TACGAGACCTATTC-1")
```



DotPlots

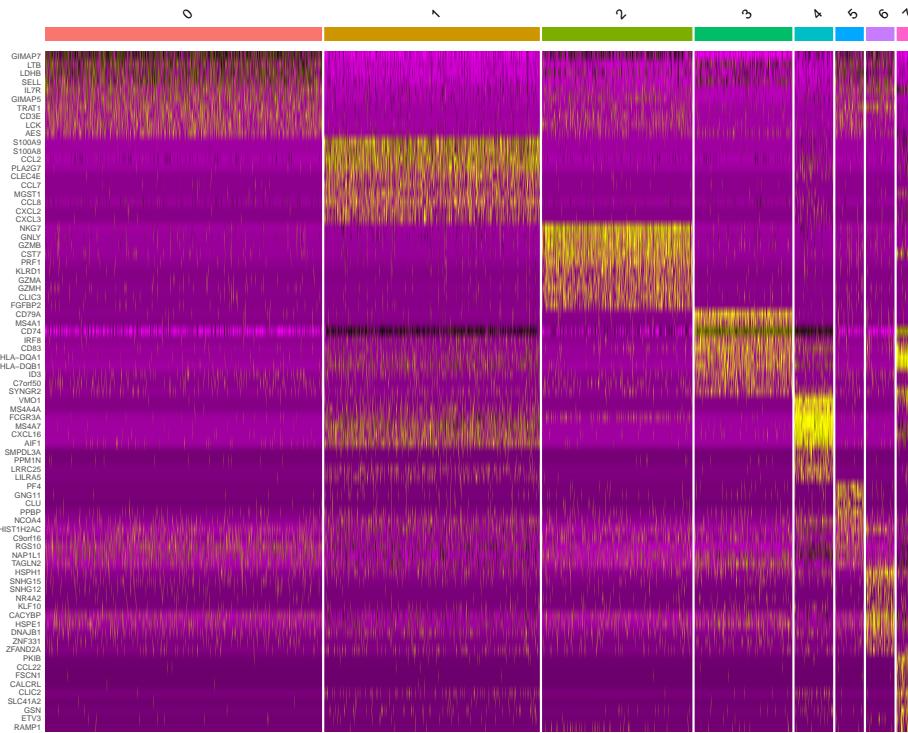
```
DotPlot(seurat_object, features = c("MS4A1", "GNLY", "CD3E", "CD14", "FCER1A", "FCGR3A"))
```

### 11.1. FINDING DIFFERENTIALLY EXPRESSED FEATURES (CLUSTER BIOMARKERS)67



`DoHeatmap()` generates an expression heatmap for given cells and features. In this case, we are plotting the top 10 markers (or all markers if less than 10) for each cluster.

```
maxcells=1500  
top10 <- seurat_object.markers %>% group_by(cluster) %>% top_n(n = 10, wt = avg_log2FC)  
DoHeatmap(subset(seurat_object, downsample = maxcells), features = top10$gene) + NoLegend()
```



## 11.2 Use makers to label or find a cluster

If you know markers for your cell types, use AddModuleScore to label them.

```
genes_markers <- list(Naive_CD4_T = c("IL7R", "CCR7"))

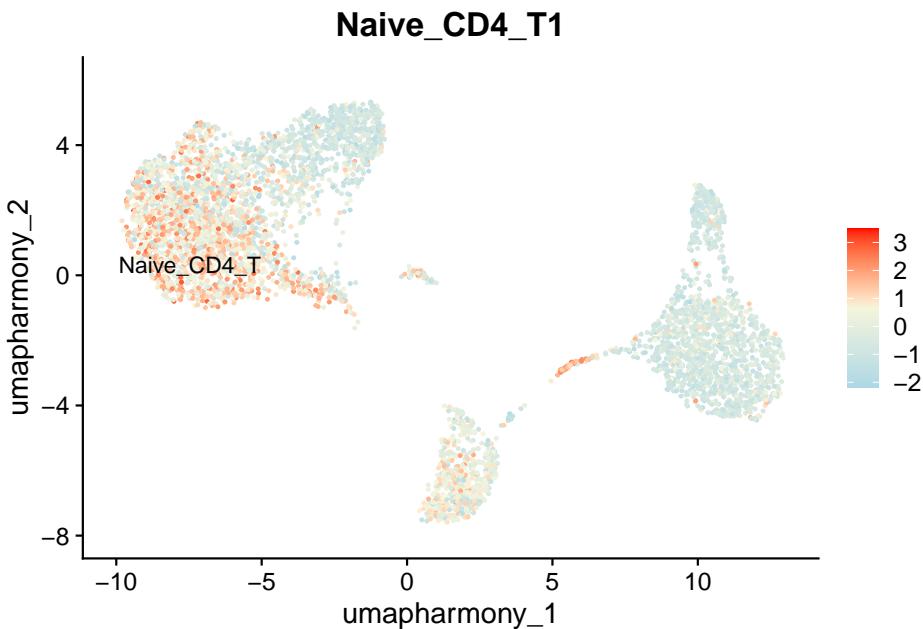
seurat_object <- AddModuleScore(object = seurat_object, features = genes_markers, ctrl
search = TRUE)

# notice the name of the cluster has a 1 at the end
names(seurat_object@meta.data)
#> [1] "orig.ident"           "nCount_RNA"
#> [3] "nFeature_RNA"        "ind"
#> [5] "stim"                 "cell"
#> [7] "multiplets"          "percent.mt"
#> [9] "RNA_snn_res.0.5"     "seurat_clusters"
#> [11] "pca_clusters"        "harmony_clusters"
#> [13] "RNA_snn_res.0.1"    "RNA_snn_res.0.2"
#> [15] "RNA_snn_res.0.3"    "RNA_snn_res.0.4"
#> [17] "RNA_snn_res.0.6"    "RNA_snn_res.0.7"
#> [19] "RNA_snn_res.0.8"    "RNA_snn_res.0.9"
```

```
#> [21] "RNA_snn_res.1"      "RNA_snn_res.1.1"
#> [23] "RNA_snn_res.1.2"    "RNA_snn_res.1.3"
#> [25] "RNA_snn_res.1.4"    "RNA_snn_res.1.5"
#> [27] "RNA_snn_res.1.6"    "RNA_snn_res.1.7"
#> [29] "RNA_snn_res.1.8"    "RNA_snn_res.1.9"
#> [31] "RNA_snn_res.2"      "Naive_CD4_T1"

# label that cell type
seurat_object$cell_label = NA
seurat_object$cell_label[seurat_object$Naive_CD4_T1 > 1] = "Naive_CD4_T"
Idents(seurat_object) = seurat_object$cell_label

# plot
# Using a custom colour scale
FeaturePlot(seurat_object, features = "Naive_CD4_T1", label = TRUE, repel = TRUE, reduction = "umap")
#> Scale for colour is already present.
#> Adding another scale for colour, which will replace the
#> existing scale.
```

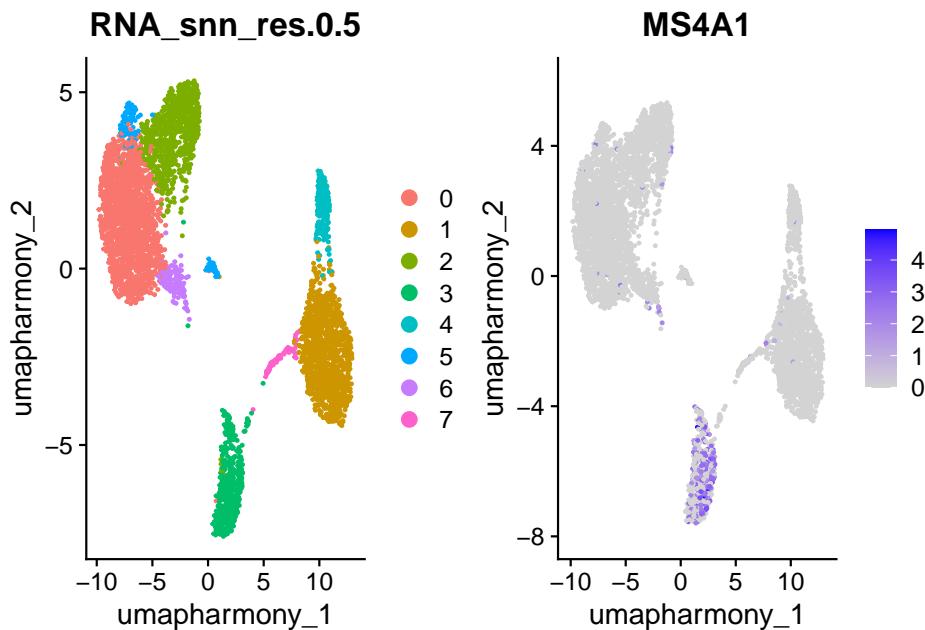


### 11.3 Assigning cell type identity to clusters

Fortunately in the case of this dataset, we can use canonical markers to easily match the unbiased clustering to known cell types:

Markers	Cell Type
IL7R, CCR7	Naive CD4+ T
CD14, LYZ	CD14+ Mono
IL7R, S100A4	Memory CD4+
MS4A1	B
CD8A	CD8+ T
FCGR3A, MS4A7	FCGR3A+ Mono
GNLY, NKG7	NK
FCER1A, CST3	DC
PPBP	Platelet

```
DimPlot(seurat_object, group.by = "RNA_snn_res.0.5", reduction = "umap_harmony") | Feature
```



### Challenge: Match cluster numbers with cell labels

Use the markers provided and the resolution 0.5 to identify the cell labels

```
Idents(seurat_object) <- seurat_object$RNA_snn_res.0.5
```

```
# This renaming list is incomplete and incorrect.
# Fix and add to the list to fill in the cell types.
new.cluster.ids <- c(
  "0" = "Naive CD4+ T",
  "2" = "CD14+ Monocytes",
  "1" = "B cells")
```

```
seurat_object <- RenameIdents(seurat_object, new.cluster.ids)
DimPlot(seurat_object, reduction = 'umap_harmony', label = TRUE, pt.size = 0.5) + NoLegend()
```

#### 11.3.0.1 save the plots

```
plot <- DimPlot(seurat_object, reduction = 'umap_harmony', label = TRUE, label.size = 4.5) + xlab(
  theme(axis.title = element_text(size = 18), legend.text = element_text(size = 18)) +
  guides(colour = guide_legend(override.aes = list(size = 10)))
ggsave(filename = "seurat_object3k_umap.jpg", height = 7, width = 12, plot = plot, quality = 50)
```

#### 11.3.0.2 save the seurat object

```
saveRDS(seurat_object, file = "seurat_object3k_final.rds")
```



## Part II

# Futher Analysis



# Chapter 12

## SingleR

```
#install.packages("BiocManager")
#BiocManager::install(c("SingleCellExperiment", "SingleR", "celldex"), ask=F)
library(SingleCellExperiment)
library(SingleR)
library(celldex)
```

In this workshop we have focused on the Seurat package. However, there is another whole ecosystem of R packages for single cell analysis within Bioconductor. We won't go into any detail on these packages in this workshop, but there is good material describing the object type online : OSCA.

For now, we'll just convert our Seurat object into an object called SingleCellExperiment. Some popular packages from Bioconductor that work with this type are Slingshot, Scran, Scater.

```
sce <- as.SingleCellExperiment(seurat_object)
sce
#> class: SingleCellExperiment
#> dim: 35635 4877
#> metadata(0):
#> assays(3): counts logcounts scaledata
#> rownames(35635): MIR1302-10 FAM138A ... MT-ND6 MT-CYB
#> rowData names(0):
#> colnames(4877): AGGGCCCTATTCC-1 GGAGACGATTCGTT-1 ...
#> ATGTTGCTAAAGC-1 GATGACACTAGCGT-1
#> colData names(34): orig.ident nCount_RNA ...
#> cell_label ident
#> reducedDimNames(4): PCA UMAP HARMONY UMAP_HARMONY
#> mainExpName: RNA
#> altExpNames(0):
```

We will now use a package called SingleR to label each cell. SingleR uses a reference data set of cell types with expression data to infer the best label for each cell. A convenient collection of cell type reference is in the `celldex` package which currently contains the follow sets:

```
ls('package:celldex')
#> [1] "BlueprintEncodeData"
#> [2] "DatabaseImmuneCellExpressionData"
#> [3] "defineTextQuery"
#> [4] "fetchLatestVersion"
#> [5] "fetchMetadata"
#> [6] "fetchReference"
#> [7] "HumanPrimaryCellAtlasData"
#> [8] "ImmGenData"
#> [9] "listReferences"
#> [10] "listVersions"
#> [11] "MonacoImmuneData"
#> [12] "MouseRNAseqData"
#> [13] "NovershternHematopoieticData"
#> [14] "saveReference"
#> [15] "searchReferences"
#> [16] "surveyReferences"
```

In this example, we'll use the `HumanPrimaryCellAtlasData` set, which contains high-level, and fine-grained label types. Lets download the reference dataset

```
# This too is a sce object,
# colData is equivalent to seurat's metadata
ref.set <- celldex::HumanPrimaryCellAtlasData()
```

The “main” labels.

```
unique(ref.set$label.main)
#> [1] "DC"                      "Smooth_muscle_cells"
#> [3] "Epithelial_cells"        "B_cell"
#> [5] "Neutrophils"             "T_cells"
#> [7] "Monocyte"                "Erythroblast"
#> [9] "BM & Prog."              "Endothelial_cells"
#> [11] "Gametocytes"            "Neurons"
#> [13] "Keratinocytes"          "HSC_-G-CSF"
#> [15] "Macrophage"              "NK_cell"
#> [17] "Embryonic_stem_cells"   "Tissue_stem_cells"
#> [19] "Chondrocytes"           "Osteoblasts"
#> [21] "BM"                     "Platelets"
#> [23] "Fibroblasts"            "iPS_cells"
#> [25] "Hepatocytes"             "MSC"
#> [27] "Neuroepithelial_cell"    "Astrocyte"
```

```
#> [29] "HSC_CD34+"           "CMP"
#> [31] "GMP"                 "MEP"
#> [33] "Myelocyte"          "Pre-B_cell_CD34-"
#> [35] "Pro-B_cell_CD34+"   "Pro-Myelocyte"
```

An example of the types of “fine” labels.

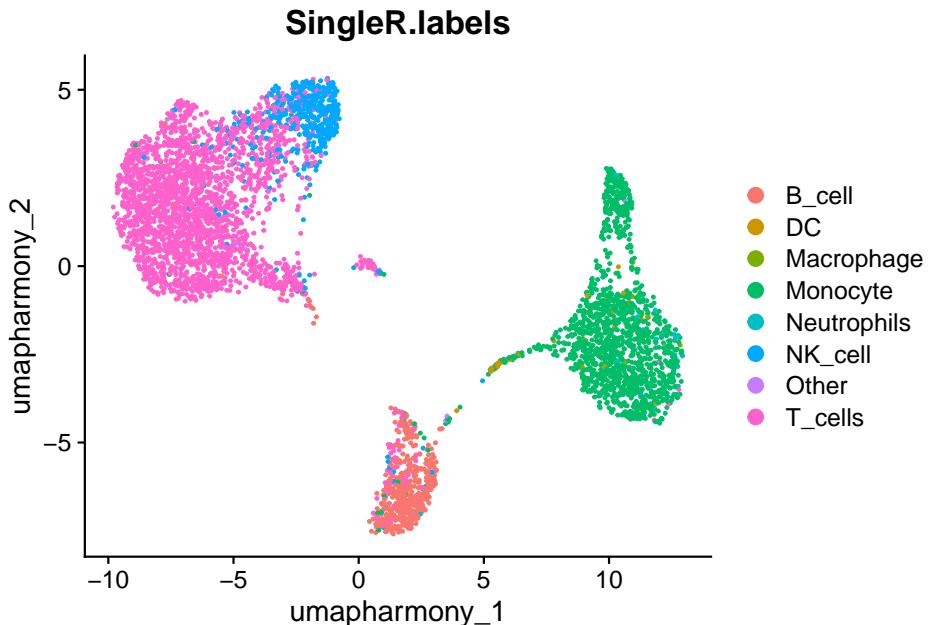
```
head(unique(ref.set$label.fine))
#> [1] "DC:monocyte-derived:immature"
#> [2] "DC:monocyte-derived:Galectin-1"
#> [3] "DC:monocyte-derived:LPS"
#> [4] "DC:monocyte-derived"
#> [5] "Smooth_muscle_cells:bronchial:vit_D"
#> [6] "Smooth_muscle_cells:bronchial"
```

Now we’ll label our cells using the SingleCellExperiment object, with the above reference set.

```
pred.cnts <- SingleR::SingleR(test = sce, ref = ref.set, labels = ref.set$label.main)
```

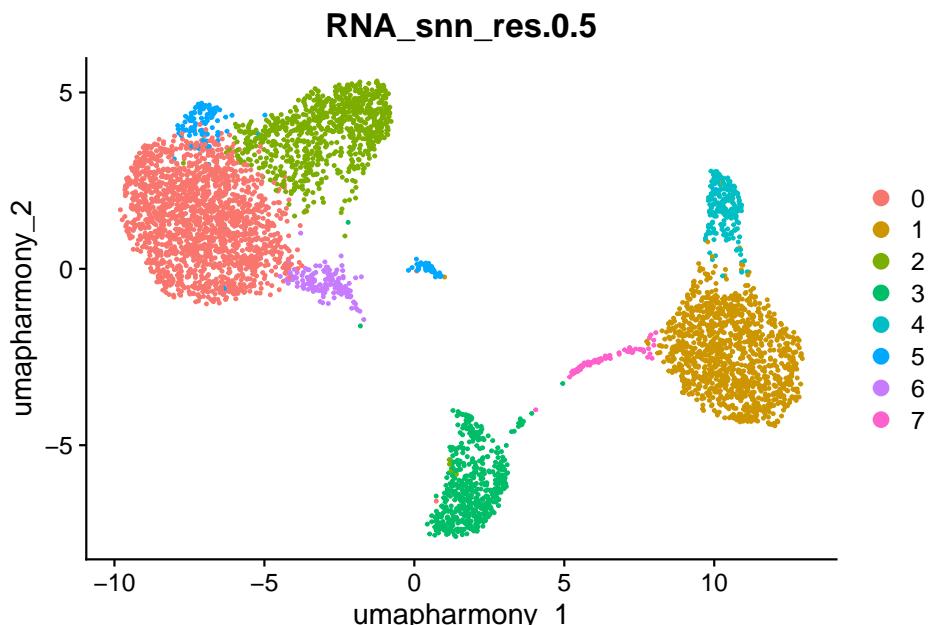
Keep any types that have more than 10 cells to the label, and put those labels back on our Seurat object and plot our on our umap.

```
lbls.keep <- table(pred.cnts$labels)>10
seurat_object$SingleR.labels <- ifelse(lbls.keep[pred.cnts$labels], pred.cnts$labels, 'Other')
DimPlot(seurat_object, reduction='umap_harmony', group.by='SingleR.labels')
```

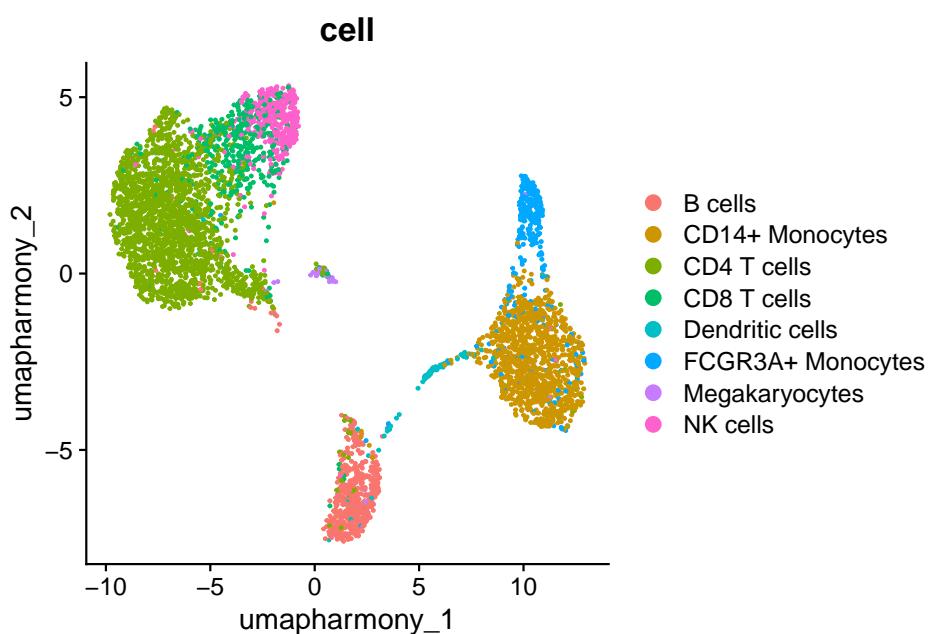


Compare cell labels by different annotation methods:

```
DimPlot(seurat_object,group.by = "RNA_snn_res.0.5",reduction = "umap_harmony")
```



```
DimPlot(seurat_object,group.by = "cell",reduction = "umap_harmony")
```



# Chapter 13

## Differential Expression

### Slides

There are many different methods for calculating differential expression between groups in scRNAseq data. There are a number of review papers worth consulting on this topic.

There is the Seurat differential expression Vignette which walks through the variety implemented in Seurat.

There is also a good discussion of using pseudobulk approaches which is worth checking out if you're planning differential expression analyses.

```
head(seurat_object@meta.data)
#>          orig.ident nCount_RNA nFeature_RNA ind
#> AGGGCGCTATTC-1 SeuratProject      2053      532 1256
#> GGAGACGATTGTT-1 SeuratProject      881       392 1256
#> CACCGTTGTCGTAG-1 SeuratProject     3130     1005 1016
#> TATCGTACACGCAT-1 SeuratProject     1042      549 1488
#> TACGAGACCTATT-1 SeuratProject     2425      777 1244
#> GTACTACTCATACG-1 SeuratProject     3951     1064 1256
#>           stim           cell multiplets
#> AGGGCGCTATTC-1 stim   CD14+ Monocytes    singlet
#> GGAGACGATTGTT-1 stim    CD4 T cells    singlet
#> CACCGTTGTCGTAG-1 ctrl FCGR3A+ Monocytes singlet
#> TATCGTACACGCAT-1 stim      B cells    singlet
#> TACGAGACCTATT-1 stim    CD4 T cells    singlet
#> GTACTACTCATACG-1 ctrl FCGR3A+ Monocytes singlet
#>           percent.mt RNA_snn_res.0.5 seurat_clusters
#> AGGGCGCTATTC-1 1.6336634            1          0
#> GGAGACGATTGTT-1 4.8809524            5         15
#> CACCGTTGTCGTAG-1 1.0655473            4         11
```

```

#> TATCGTACACGCAT-1 3.0662710      3      5
#> TACGAGACCTATTTC-1 1.0837849      0      9
#> GTACTACTCATACG-1 0.7137395      4     11
#>          pca_clusters harmony_clusters
#> AGGGCGCTATTTC-1      3      1
#> GGAGACGATTCTGTT-1      0      5
#> CACCGTTGTCGTAG-1      9      4
#> TATCGTACACGCAT-1      6      3
#> TACGAGACCTATTTC-1      0      0
#> GTACTACTCATACG-1      9      4
#>          RNA_snn_res.0.1 RNA_snn_res.0.2
#> AGGGCGCTATTTC-1      1      1
#> GGAGACGATTCTGTT-1      0      0
#> CACCGTTGTCGTAG-1      3      4
#> TATCGTACACGCAT-1      2      3
#> TACGAGACCTATTTC-1      0      0
#> GTACTACTCATACG-1      3      4
#>          RNA_snn_res.0.3 RNA_snn_res.0.4
#> AGGGCGCTATTTC-1      1      1
#> GGAGACGATTCTGTT-1      0      0
#> CACCGTTGTCGTAG-1      4      4
#> TATCGTACACGCAT-1      3      3
#> TACGAGACCTATTTC-1      0      0
#> GTACTACTCATACG-1      4      4
#>          RNA_snn_res.0.6 RNA_snn_res.0.7
#> AGGGCGCTATTTC-1      1      1
#> GGAGACGATTCTGTT-1      7      6
#> CACCGTTGTCGTAG-1      5      5
#> TATCGTACACGCAT-1      3      3
#> TACGAGACCTATTTC-1      0      0
#> GTACTACTCATACG-1      5      5
#>          RNA_snn_res.0.8 RNA_snn_res.0.9
#> AGGGCGCTATTTC-1      1      0
#> GGAGACGATTCTGTT-1      7      8
#> CACCGTTGTCGTAG-1      6      7
#> TATCGTACACGCAT-1      4      3
#> TACGAGACCTATTTC-1      0      1
#> GTACTACTCATACG-1      6      7
#>          RNA_snn_res.1 RNA_snn_res.1.1
#> AGGGCGCTATTTC-1      1      0
#> GGAGACGATTCTGTT-1      9      9
#> CACCGTTGTCGTAG-1      7      7
#> TATCGTACACGCAT-1      3      3
#> TACGAGACCTATTTC-1      0      1
#> GTACTACTCATACG-1      7      7

```

```

#>          RNA_snn_res.1.2 RNA_snn_res.1.3
#> AGGGCGCTATTCC-1      2      0
#> GGAGACGATTCTGTT-1   9     12
#> CACCGTTGTCGTAG-1    7      9
#> TATCGTACACGCAT-1    3      2
#> TACGAGACCTATTCC-1   0      8
#> GTACTACTCATACG-1    7      9
#>          RNA_snn_res.1.4 RNA_snn_res.1.5
#> AGGGCGCTATTCC-1      0      0
#> GGAGACGATTCTGTT-1   14     14
#> CACCGTTGTCGTAG-1    10     11
#> TATCGTACACGCAT-1    2      2
#> TACGAGACCTATTCC-1   8      9
#> GTACTACTCATACG-1    10     11
#>          RNA_snn_res.1.6 RNA_snn_res.1.7
#> AGGGCGCTATTCC-1      2      0
#> GGAGACGATTCTGTT-1   15     14
#> CACCGTTGTCGTAG-1    12     11
#> TATCGTACACGCAT-1    0      2
#> TACGAGACCTATTCC-1   9      9
#> GTACTACTCATACG-1    12     11
#>          RNA_snn_res.1.8 RNA_snn_res.1.9
#> AGGGCGCTATTCC-1      0      4
#> GGAGACGATTCTGTT-1   14     16
#> CACCGTTGTCGTAG-1    11     12
#> TATCGTACACGCAT-1    2      0
#> TACGAGACCTATTCC-1   10     6
#> GTACTACTCATACG-1    11     12
#>          RNA_snn_res.2 Naive_CD4_T1 cell_label
#> AGGGCGCTATTCC-1      0   -0.3540023 <NA>
#> GGAGACGATTCTGTT-1   15   2.1376158 Naive_CD4_T
#> CACCGTTGTCGTAG-1    11  -1.1487836 <NA>
#> TATCGTACACGCAT-1    5   0.5557941 <NA>
#> TACGAGACCTATTCC-1   9   1.4250065 Naive_CD4_T
#> GTACTACTCATACG-1    11  -0.2082793 <NA>
#>          SingleR.labels
#> AGGGCGCTATTCC-1      Monocyte
#> GGAGACGATTCTGTT-1    T_cells
#> CACCGTTGTCGTAG-1     Monocyte
#> TATCGTACACGCAT-1     Neutrophils
#> TACGAGACCTATTCC-1    T_cells
#> GTACTACTCATACG-1     Monocyte

```

How cells from each condition do we have?

```
table(seurat_object$stim)
#>
#> ctrl stim
#> 2378 2499
```

How many cells per individuals per group?

```
table(seurat_object$ind, seurat_object$stim)
#>
#>          ctrl stim
#> 101     175 226
#> 107     115 106
#> 1015    503 489
#> 1016    335 348
#> 1039    87 100
#> 1244    373 311
#> 1256    384 385
#> 1488    406 534
```

And for each sample, how many of each cell type has been classified?

```
table(paste(seurat_object$ind, seurat_object$stim), seurat_object$cell)
#>
#>          B cells CD14+ Monocytes CD4 T cells CD8 T cells
#> 101 ctrl      23        47       60       15
#> 101 stim      30        52       83       17
#> 1015 ctrl     79        149      142       44
#> 1015 stim     68        149      148       21
#> 1016 ctrl     21        71       82      107
#> 1016 stim     29        65       65      109
#> 1039 ctrl      7        32       36        6
#> 1039 stim      6        26       48        5
#> 107 ctrl       9        50       32        6
#> 107 stim       9        35       43        1
#> 1244 ctrl     23        85      202        8
#> 1244 stim     18        58      191        4
#> 1256 ctrl     32        80      176       26
#> 1256 stim     42        70      194       25
#> 1488 ctrl     36        59      246       13
#> 1488 stim     59        59      319       15
#>
#>          Dendritic cells FCGR3A+ Monocytes
#> 101 ctrl        4        11
#> 101 stim        6        23
#> 1015 ctrl       3        49
#> 1015 stim      17        43
```

```

#> 1016 ctrl      4      21
#> 1016 stim     2      32
#> 1039 ctrl      1      3
#> 1039 stim     1      6
#> 107 ctrl       3      12
#> 107 stim       2      5
#> 1244 ctrl      8      19
#> 1244 stim     6      4
#> 1256 ctrl      6      20
#> 1256 stim     3      11
#> 1488 ctrl      8      25
#> 1488 stim     12      28
#>
#>          Megakaryocytes NK cells
#> 101 ctrl       4      11
#> 101 stim       1      14
#> 1015 ctrl      5      32
#> 1015 stim      5      38
#> 1016 ctrl      3      26
#> 1016 stim      1      45
#> 1039 ctrl      1      1
#> 1039 stim      3      5
#> 107 ctrl       0      3
#> 107 stim       0      11
#> 1244 ctrl      3      25
#> 1244 stim      4      26
#> 1256 ctrl      1      43
#> 1256 stim      7      33
#> 1488 ctrl      4      15
#> 1488 stim      6      36

```

## 13.1 Prefiltering

### Why do we need to do this?

If expression is below a certain level, it will be almost impossible to see any differential expression.

When doing differential expression, you generally ignore genes with low expression. In single cell datasets, there are many genes like this. Filtering here to make our dataset smaller so it runs quicker, and there is less aggressive correction for multiple hypotheses.

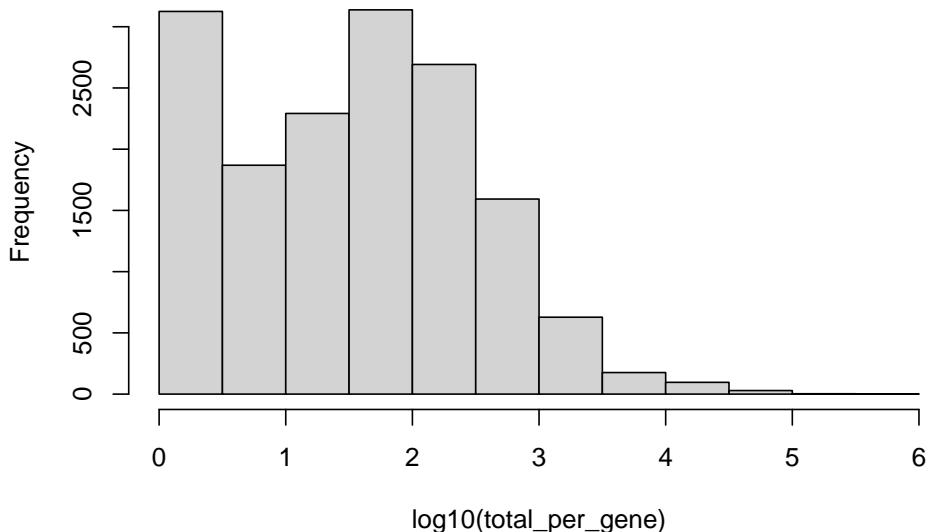
How many genes before filtering?

```
seurat_object
#> An object of class Seurat
#> 35635 features across 4877 samples within 1 assay
#> Active assay: RNA (35635 features, 2000 variable features)
#> 3 layers present: counts, data, scale.data
#> 4 dimensional reductions calculated: pca, umap, harmony, umap_harmony
```

How many copies of each gene are there?

```
total_per_gene <- rowSums(GetAssayData(seurat_object, assay='RNA', slot='counts'))
hist(log10(total_per_gene))
```

**Histogram of log10(total\_per\_gene)**



Lets keep only those genes with at least 50 copies across the entire experiment.

```
seurat_object <- seurat_object[total_per_gene >= 50, ]
```

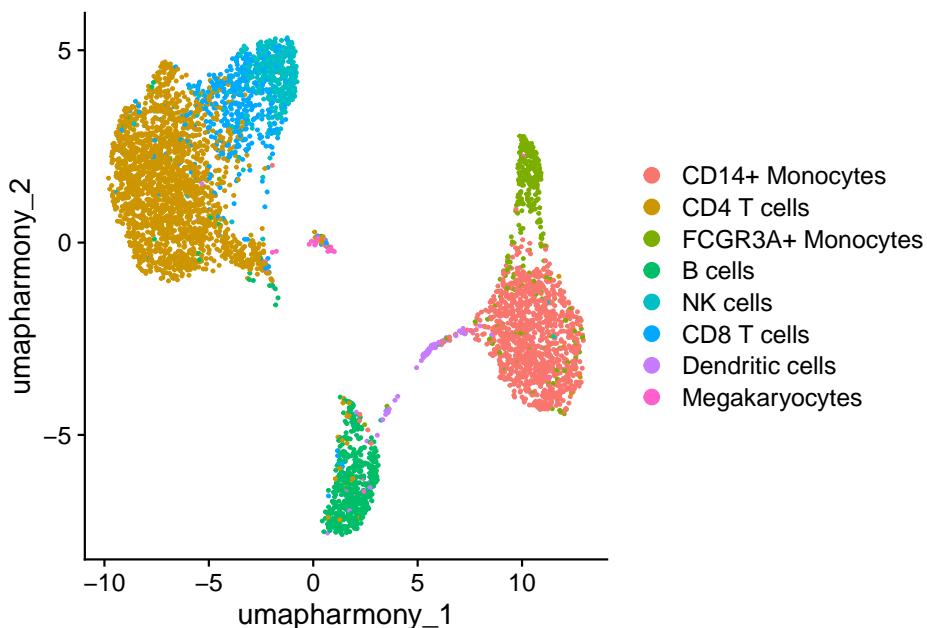
How many genes after filtering?

```
seurat_object
#> An object of class Seurat
#> 7200 features across 4877 samples within 1 assay
#> Active assay: RNA (7200 features, 1236 variable features)
#> 3 layers present: counts, data, scale.data
#> 4 dimensional reductions calculated: pca, umap, harmony, umap_harmony
```

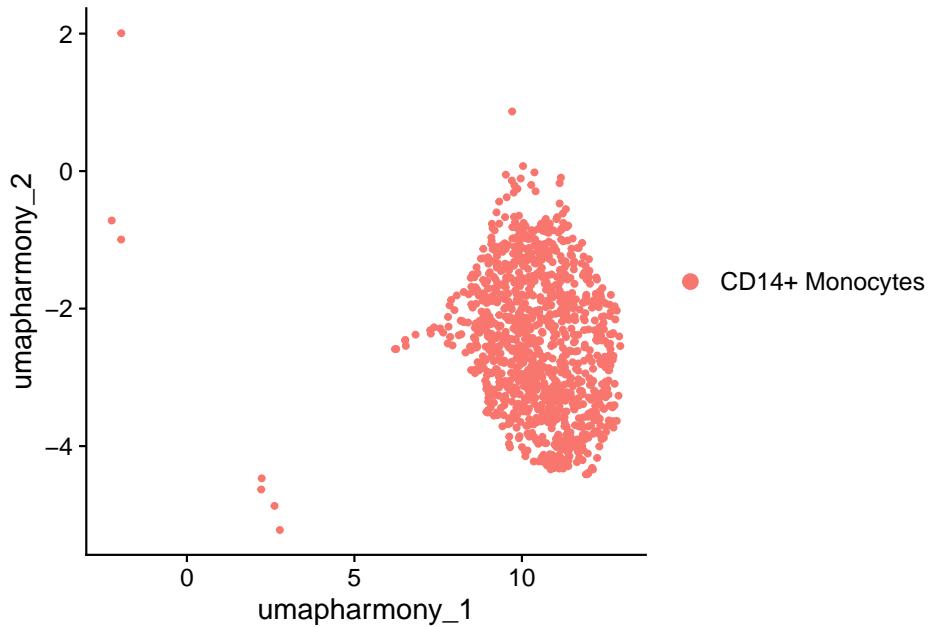
We might like to see the effect of IFN-beta stimulation on each cell type individually. For the purposes of this workshop, just going to test one cell type; CD14+ Monocytes

An easy way is to subset the object.

```
# Set idents to 'cell' column.  
Idents(seurat_object) <- seurat_object$cell  
DimPlot(seurat_object, reduction = "umap_harmony")
```



```
seurat_object_celltype <- seurat_object[, seurat_object$cell == "CD14+ Monocytes" ]  
DimPlot(seurat_object_celltype, reduction = "umap_harmony")
```



## 13.2 Default Wilcox test

To run this test, we change the Idents to the factor(column) we want to test. In this case, that's 'stim'.

```
# Change Ident to Condition
Idents(seurat_object_celltype) <- seurat_object_celltype$stim

# default, wilcox test
de_result_wilcox <- FindMarkers(seurat_object_celltype,
  ident.1 = 'stim',
  ident.2 = 'ctrl',
  logfc.threshold = 0, # Give me ALL results
  min.pct = 0
)

# Add average expression for plotting
de_result_wilcox$AveExpr<- rowMeans(seurat_object_celltype[rownames(de_result_wilcox),])
```

Look at the top differentially expressed genes.

```
head(de_result_wilcox)
#>          p_val avg_log2FC pct.1 pct.2      p_val_adj
#> RSAD2    8.471603e-196   6.797508 0.979 0.044 6.099554e-192
#> TNFSF10  3.399691e-195   6.559853 0.982 0.056 2.447777e-191
#> CXCL10   5.809758e-195   8.034949 0.975 0.038 4.183025e-191
```

```

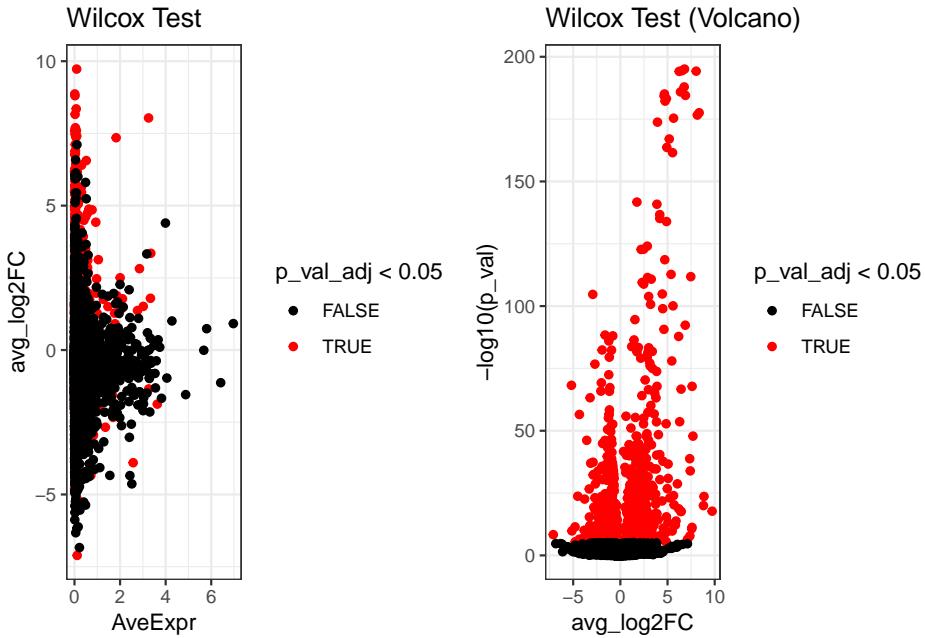
#> IFIT3    7.214495e-195   6.234703 0.982 0.051 5.194437e-191
#> IFIT1    1.224398e-188   6.764699 0.953 0.026 8.815662e-185
#> ISG15    1.228738e-186   6.377243 1.000 0.323 8.846914e-183
#>          AveExpr
#> RSAD2    0.04196286
#> TNFSF10  0.51655832
#> CXCL10   3.25184139
#> IFIT3    0.02032398
#> IFIT1    0.01254584
#> ISG15    0.14165156

p1 <- ggplot(de_result_wilcox, aes(x=AveExpr, y=avg_log2FC, col=p_val_adj < 0.05)) +
  geom_point() +
  scale_colour_manual(values=c('TRUE'="red", 'FALSE'="black")) +
  theme_bw() +
  ggtitle("Wilcox Test")

p2 <- ggplot(de_result_wilcox, aes(x=avg_log2FC, y=-log10(p_val), col=p_val_adj < 0.05)) +
  geom_point() +
  scale_colour_manual(values=c('TRUE'="red", 'FALSE'="black")) +
  theme_bw() +
  ggtitle("Wilcox Test (Volcano)")

p1 + p2

```



### 13.3 Seurat Negative binomial

Negative binomial test is run almost the same way - just need to specify it under 'test.use'

```
# Change Ident to Condition
Idents(seurat_object_celltype) <- seurat_object_celltype$stim

# default, wilcox test
de_result_negbinom <- FindMarkers(seurat_object_celltype,
                                    test.use="negbinom", # Choose a different test.
                                    ident.1 = 'stim',
                                    ident.2 = 'ctrl',
                                    logfc.threshold = 0, # Give me ALL results
                                    min.pct = 0
)

# Add average expression for plotting
de_result_negbinom$AveExpr<- rowMeans(seurat_object_celltype[rownames(de_result_negbinom),]
```

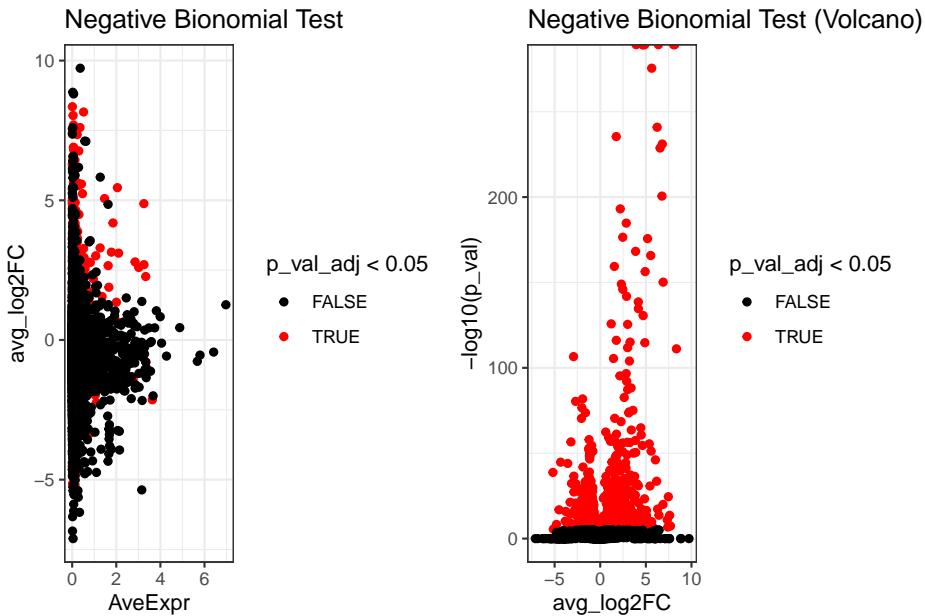
Look at the top differentially expressed genes.

```
head(de_result_negbinom)
#>      p_val avg_log2FC pct.1 pct.2 p_val_adj     AveExpr
#> CXCL10      0   8.034949 0.975 0.038      0 0.04196286
#> CCL8        0   8.160755 0.912 0.023      0 0.51655832
#> LY6E        0   4.880093 0.979 0.134      0 3.25184139
#> APOBEC3A    0   4.745713 0.996 0.262      0 0.02032398
#> IFITM3      0   4.675898 0.998 0.271      0 0.01254584
#> IFI6        0   3.935122 0.982 0.257      0 0.14165156

p1 <- ggplot(de_result_negbinom, aes(x=AveExpr, y=avg_log2FC, col=p_val_adj < 0.05)) +
  geom_point() +
  scale_colour_manual(values=c('TRUE'="red", 'FALSE'="black")) +
  theme_bw() +
  ggtitle("Negative Bionomial Test")

p2 <- ggplot(de_result_negbinom, aes(x=avg_log2FC, y=-log10(p_val), col=p_val_adj < 0.05)) +
  geom_point() +
  scale_colour_manual(values=c('TRUE'="red", 'FALSE'="black")) +
  theme_bw() +
  ggtitle("Negative Bionomial Test (Volcano)")

p1 + p2
```



## 13.4 Pseudobulk

Pseudobulk analysis is an option where you have biological replicates. It is essentially pooling the individual cell counts and treating your experiment like a bulk RNAseq.

First, you need to build a pseudobulk matrix - the `AggregateExpression()` function can do this, once you set the ‘Idents’ of your seurat object to your grouping factor (here, that’s a combination of individual+treatment called ‘sample’, instead of the ‘stim’ treatment column).

```
# Tools for bulk differential expression
library(limma)
#>
#> Attaching package: 'limma'
#> The following object is masked from 'package:BiocGenerics':
#>
#>     plotMA
library(edgeR)
#>
#> Attaching package: 'edgeR'
#> The following object is masked from 'package:SingleCellExperiment':
#>
#>     cpm
```

```
# Change idents to ind for grouping.
seurat_object_celltype$sample <- factor(paste(seurat_object_celltype$stim, seurat_object_celltype$sample))
Idents(seurat_object_celltype) <- seurat_object_celltype$sample

# Then pool together counts in those groups
# AggregateExpression returns a list of matrices - one for each assay requested (even
pseudobulk_matrix_list <- AggregateExpression(seurat_object_celltype, slot = 'counts')
## Names of identity class contain underscores ('_'), replacing with dashes ('-')
## This message is displayed once every 8 hours.
pseudobulk_matrix      <- pseudobulk_matrix_list[['RNA']]
pseudobulk_matrix[1:5,1:4]
## 5 x 4 sparse Matrix of class "dgCMatrix"
##>          ctrl-101 ctrl-1015 ctrl-1016 ctrl-1039
##> NOC2L           2       7       .       .
##> HES4            .       3       2       1
##> ISG15          31     185     234     41
##> TNFRSF18        .       3       4       2
##> TNFRSF4          .       2       .       .
```

Now it looks like a bulk RNAseq experiment, so treat it like one.

We can use the popular `limma` package for differential expression. Here is one tutorial, and the hefty reference manual is hosted by bioconductor.

In brief, this code below constructs a linear model for this experiment that accounts for the variation in individuals and treatment. It then tests for differential expression between 'stim' and 'ctrl' groups.

```
dge <- DGEList(pseudobulk_matrix)
dge <- calcNormFactors(dge)

## Remove _ or - and everything after it - yeilds stim group
stim <- gsub("[-_].*", "", colnames(pseudobulk_matrix))

## Removing everything before the _ or - for the individual, then converting those numbers
ind  <- as.character(gsub(".*[-_]", "", colnames(pseudobulk_matrix)))

design <- model.matrix(~0 + stim + ind)
vm   <- voom(dge, design = design, plot = FALSE)
fit <- lmFit(vm, design = design)

contrasts <- makeContrasts(stimstim - stimctrl, levels=coef(fit))
fit <- contrasts.fit(fit, contrasts)
fit <- eBayes(fit)

de_result_pseudobulk <- topTable(fit, n = Inf, adjust.method = "BH")
de_result_pseudobulk <- arrange(de_result_pseudobulk , adj.P.Val)
```

Look at the significantly differentially expressed genes:

```

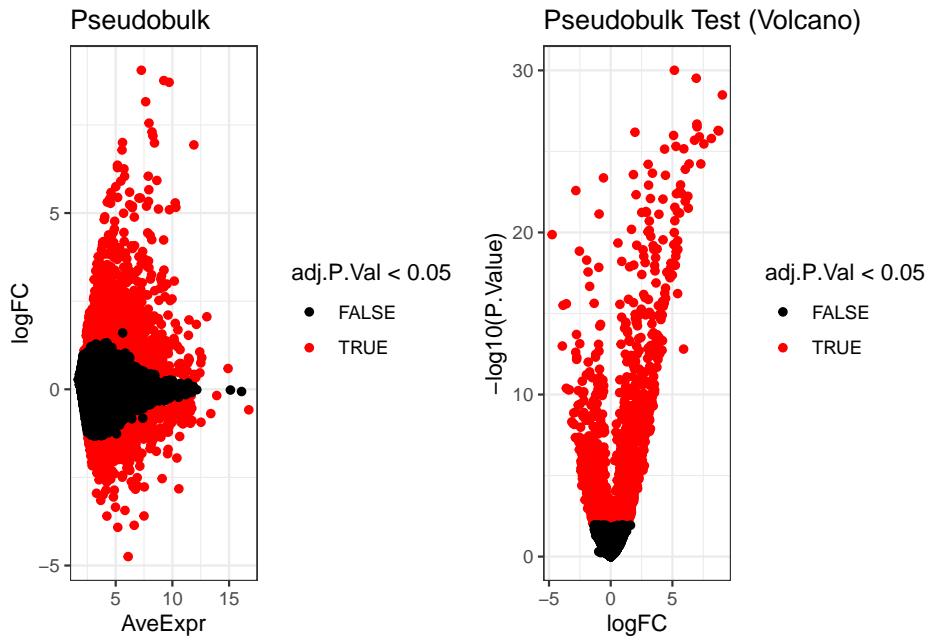
head(de_result_pseudobulk)
#>      logFC    AveExpr      t     P.Value
#> ISG20  5.161336 10.313033 35.72209 9.838370e-31
#> ISG15  6.932960 11.897783 34.63778 3.049704e-30
#> CXCL11 9.050139  7.263356 32.46019 3.282727e-29
#> IFIT3  6.987552  8.423467 28.97126 2.052443e-27
#> HERC5  6.999781  5.604886 28.65574 3.050588e-27
#> CXCL10 8.708460  9.715253 28.25406 5.081902e-27
#>          adj.P.Val      B
#> ISG20  7.083626e-27 60.13277
#> ISG15  1.097894e-26 59.02083
#> CXCL11 7.878545e-26 55.45779
#> IFIT3  3.694397e-24 52.01990
#> HERC5  4.392847e-24 51.18902
#> CXCL10 5.629458e-24 51.32019

p1 <- ggplot(de_result_pseudobulk, aes(x=AveExpr, y=logFC, col=adj.P.Val < 0.05)) +
  geom_point() +
  scale_colour_manual(values=c('TRUE'='red', 'FALSE'='black')) +
  theme_bw() +
  ggtitle("Pseudobulk")

p2 <- ggplot(de_result_pseudobulk, aes(x=logFC, y=-log10(P.Value), col=adj.P.Val < 0.05)) +
  geom_point() +
  scale_colour_manual(values=c('TRUE'='red', 'FALSE'='black')) +
  theme_bw() +
  ggtitle("Pseudobulk Test (Volcano)")

p1 + p2

```



### Discussion

These methods give different results. How would you decide which to use? How could you check an individual gene?

# Chapter 14

## Cell cycle Assignment

In some datasets, the phase of cell cycle that a cell is in (G1/G2M/S) can account for a lot of the observed transcriptomic variation. There may be clustering by phase, or separation in the UMAP by phase.

Seurat provides a simple method for assigning cell cycle state to each cell. Other methods are available.

More information about assigning cell cycle states to cells is in the cell cycle vignette

```
# A list of cell cycle markers, from Tirosh et al, 2015, is loaded with Seurat. We can
# segregate this list into markers of G2/M phase and markers of S phase
s.genes <- cc.genes$s.genes
g2m.genes <- cc.genes$g2m.genes

# Use those lists with the cell cycle scoring function in Seurat.
seurat_object <- CellCycleScoring(seurat_object, s.features = s.genes, g2m.features = g2m.genes)
#> Warning: The following features are not present in the
#> object: TYMS, MCM2, MCM4, UNG, GINS2, CDCA7, DTL, UHRF1,
#> MLF1IP, HELLS, RAD51AP1, GMNN, WDR76, CCNE2, ATAD2, RAD51,
#> RRM2, CDC45, CDC6, EXO1, DSCC1, BLM, CASP8AP2, CLSPN,
#> POLA1, CHAF1B, BRIP1, E2F8, not searching for symbol
#> synonyms
#> Warning: The following features are not present in the
#> object: CDK1, NUSAP1, UBE2C, BIRC5, TPX2, TOP2A, NDC80,
#> NUF2, MKI67, FAM64A, CCNB2, CKAP2L, AURKB, BUB1, KIF11,
#> GTSE1, HJURP, CDCA3, CDC20, TTK, CDC25C, KIF2C, NCAPD2,
#> DLGAP5, CDCA2, CDCA8, ECT2, KIF23, HMMR, AURKA, PSRC1,
#> ANLN, CENPE, NEK2, GAS2L3, CENPA, not searching for symbol
#> synonyms
```

Which adds S.Score, G2M.Score and Phase calls to the metadata.

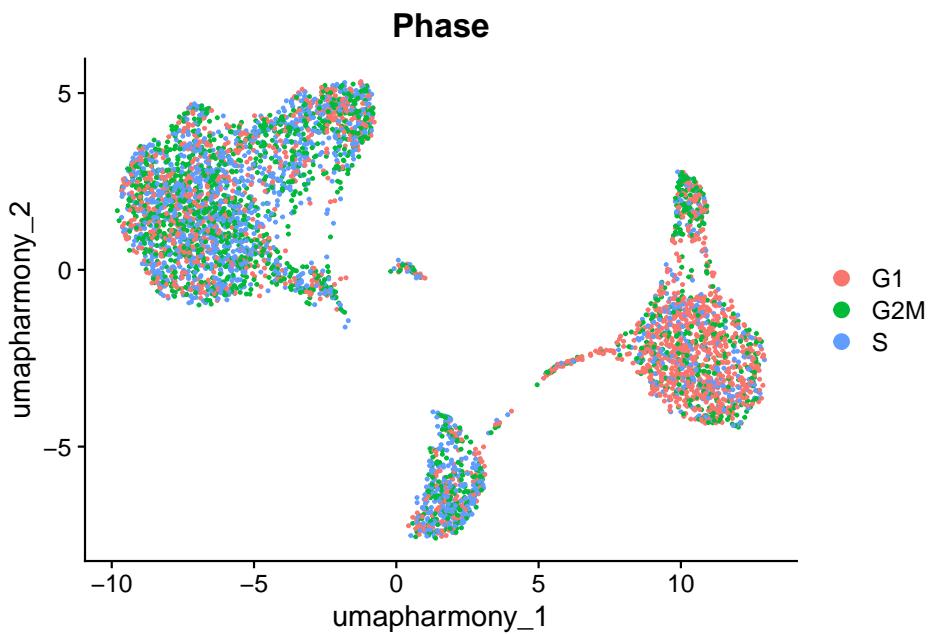
```
head(seurat_object@meta.data)
#> orig.ident nCount_RNA nFeature_RNA ind
#> AGGGCGCTATTTCC-1 SeuratProject      2053      532 1256
#> GGAGACGATTCTCGTT-1 SeuratProject     881       392 1256
#> CACCGTTGTCGTAG-1 SeuratProject      3130      1005 1016
#> TATCGTACACGCAT-1 SeuratProject      1042       549 1488
#> TACGAGACCTATTTC-1 SeuratProject     2425       777 1244
#> GTACTACTCATACG-1 SeuratProject      3951      1064 1256
#>          stim      cell multiplets
#> AGGGCGCTATTTCC-1 stim    CD14+ Monocytes   singlet
#> GGAGACGATTCTCGTT-1 stim    CD4 T cells     singlet
#> CACCGTTGTCGTAG-1 ctrl FCGR3A+ Monocytes   singlet
#> TATCGTACACGCAT-1 stim    B cells        singlet
#> TACGAGACCTATTTC-1 stim    CD4 T cells     singlet
#> GTACTACTCATACG-1 ctrl FCGR3A+ Monocytes   singlet
#>          percent.mt RNA_snn_res.0.5 seurat_clusters
#> AGGGCGCTATTTCC-1 1.6336634           1         0
#> GGAGACGATTCTCGTT-1 4.8809524           5         15
#> CACCGTTGTCGTAG-1 1.0655473           4         11
#> TATCGTACACGCAT-1 3.0662710           3         5
#> TACGAGACCTATTTC-1 1.0837849           0         9
#> GTACTACTCATACG-1 0.7137395           4         11
#>          pca_clusters harmony_clusters
#> AGGGCGCTATTTCC-1      3         1
#> GGAGACGATTCTCGTT-1     0         5
#> CACCGTTGTCGTAG-1      9         4
#> TATCGTACACGCAT-1      6         3
#> TACGAGACCTATTTC-1     0         0
#> GTACTACTCATACG-1      9         4
#>          RNA_snn_res.0.1 RNA_snn_res.0.2
#> AGGGCGCTATTTCC-1      1         1
#> GGAGACGATTCTCGTT-1     0         0
#> CACCGTTGTCGTAG-1      3         4
#> TATCGTACACGCAT-1      2         3
#> TACGAGACCTATTTC-1     0         0
#> GTACTACTCATACG-1      3         4
#>          RNA_snn_res.0.3 RNA_snn_res.0.4
#> AGGGCGCTATTTCC-1      1         1
#> GGAGACGATTCTCGTT-1     0         0
#> CACCGTTGTCGTAG-1      4         4
#> TATCGTACACGCAT-1      3         3
#> TACGAGACCTATTTC-1     0         0
#> GTACTACTCATACG-1      4         4
#>          RNA_snn_res.0.6 RNA_snn_res.0.7
```

#> AGGGCGCTATTC-1	1	1
#> GGAGACGATTGTT-1	7	6
#> CACCGTTGTCGTAG-1	5	5
#> TATCGTACACGCAT-1	3	3
#> TACGAGACCTATT-1	0	0
#> GTACTACTCATACG-1	5	5
	<i>RNA_snn_res.0.8 RNA_snn_res.0.9</i>	
#> AGGGCGCTATTC-1	1	0
#> GGAGACGATTGTT-1	7	8
#> CACCGTTGTCGTAG-1	6	7
#> TATCGTACACGCAT-1	4	3
#> TACGAGACCTATT-1	0	1
#> GTACTACTCATACG-1	6	7
	<i>RNA_snn_res.1 RNA_snn_res.1.1</i>	
#> AGGGCGCTATTC-1	1	0
#> GGAGACGATTGTT-1	9	9
#> CACCGTTGTCGTAG-1	7	7
#> TATCGTACACGCAT-1	3	3
#> TACGAGACCTATT-1	0	1
#> GTACTACTCATACG-1	7	7
	<i>RNA_snn_res.1.2 RNA_snn_res.1.3</i>	
#> AGGGCGCTATTC-1	2	0
#> GGAGACGATTGTT-1	9	12
#> CACCGTTGTCGTAG-1	7	9
#> TATCGTACACGCAT-1	3	2
#> TACGAGACCTATT-1	0	8
#> GTACTACTCATACG-1	7	9
	<i>RNA_snn_res.1.4 RNA_snn_res.1.5</i>	
#> AGGGCGCTATTC-1	0	0
#> GGAGACGATTGTT-1	14	14
#> CACCGTTGTCGTAG-1	10	11
#> TATCGTACACGCAT-1	2	2
#> TACGAGACCTATT-1	8	9
#> GTACTACTCATACG-1	10	11
	<i>RNA_snn_res.1.6 RNA_snn_res.1.7</i>	
#> AGGGCGCTATTC-1	2	0
#> GGAGACGATTGTT-1	15	14
#> CACCGTTGTCGTAG-1	12	11
#> TATCGTACACGCAT-1	0	2
#> TACGAGACCTATT-1	9	9
#> GTACTACTCATACG-1	12	11
	<i>RNA_snn_res.1.8 RNA_snn_res.1.9</i>	
#> AGGGCGCTATTC-1	0	4
#> GGAGACGATTGTT-1	14	16
#> CACCGTTGTCGTAG-1	11	12

```
#> TATCGTACACGCAT-1          2          0
#> TACGAGACCTATTTC-1        10         6
#> GTACTACTCATACG-1         11         12
#> RNA_snn_res.2 Naive_CD4_T1  cell_label
#> AGGGCGCTATTTCC-1          0   -0.3540023      <NA>
#> GGAGACGATTCTGTT-1         15   2.1376158 Naive_CD4_T
#> CACCGTTGTCGTAG-1          11  -1.1487836      <NA>
#> TATCGTACACGCAT-1          5    0.5557941      <NA>
#> TACGAGACCTATTTC-1         9    1.4250065 Naive_CD4_T
#> GTACTACTCATACG-1          11  -0.2082793      <NA>
#> SingleR.labels      S.Score  G2M.Score
#> AGGGCGCTATTTCC-1  Monocyte -0.06808986 -0.08878508
#> GGAGACGATTCTGTT-1  T_cells  0.09975063 -0.06812238
#> CACCGTTGTCGTAG-1  Monocyte -0.12770457  0.16052598
#> TATCGTACACGCAT-1  Neutrophils 0.08685571 -0.08316351
#> TACGAGACCTATTTC-1  T_cells  -0.02572695  0.02074759
#> GTACTACTCATACG-1  Monocyte -0.07294007 -0.01517046
#> Phase
#> AGGGCGCTATTTCC-1    G1
#> GGAGACGATTCTGTT-1    S
#> CACCGTTGTCGTAG-1    G2M
#> TATCGTACACGCAT-1    S
#> TACGAGACCTATTTC-1   G2M
#> GTACTACTCATACG-1    G1
```

We can then check the cell phase on the UMAP. In this dataset, phase isn't driving the clustering, and would not require any further handling.

```
DimPlot(seurat_object, reduction = 'umap_harmony', group.by = "Phase")
```



Where a bias *is* present, your course of action depends on the task at hand. It might involve ‘regressing out’ the cell cycle variation when scaling data `ScaleData(kang, vars.to.regress="Phase")`, omitting cell-cycle dominated clusters, or just accounting for it in your differential expression calculations.

If you are working with non-human data, you will need to convert these gene lists, or find new cell cycle associated genes in your species.



## **Part III**

# **Other Resources**



# Chapter 15

## Resources

Useful resources for next steps.

### 15.0.1 Suggested Further Reading Material

- Orchestrating Single Cell Analysis with Bioconductor - this book teaches single cell analysis with the bioconductor ecosystem of packages rather than Seurat. Regardless of your preference for Bioconductor or Seurat, it provides an excellent grounding and further depth and rationale behind each step of a single cell analysis.
- Seurat tutorials for gene expression, spatial & multimodal analysis
- Getting started with Signac - the sibling package to Seurat for scATAC analysis
- Monocle documentationn for trajectories
- Cell Annotation with SingleR
- VDJ analysis with Immcantation

### 15.0.2 Useful links arising from the discussion during the previous workshop

- 10x Genomics link to ribosomal protein expression
- 10x Genomics link to mitochondrial gene expression
- scRNA Tools, catalogue of tools for scRNA Seq analysis

#### 15.0.2.1 Data interpretation

- Interactive website explaining UMAP and comparision to t-SNE.
- OSCA, dimensionality reduction interpretation
- A simple description of what PCA and UMAP do, with a 3D example.

### 15.0.2.2 Data tools and visualisation

- scTransform Vignette
- Link to the workflow library
- iSEE Bioconductor library, interactive explorer
- ShinyCell makes interactive Shiny app from Seurat output
- iCellR interactive data explorer
- Diffusion maps for single cell instead of umaps \*Projections of a high-dimensional dataset with an animated scatter-plot

### 15.0.2.3 Papers

- Doublet cell detection method benchmarking paper.
- From Louvain to Leiden: guaranteeing well-connected communities

### 15.0.2.4 Reference data and databases

- Gene tissue expression database
- ImmGen Database and Explorer
- Single Cell Study Portal from The Broad
- Common ref data for cell indexing
- Azimuth is a Seurat-friendly reference-based annotation tool
- Celaref, cell reference annotation tool

## 15.1 Help and further Resources

### Seurat Vignettes

<https://satijalab.org/seurat/index.html>

There are a good many Seurat vignettes for different aspects of the Seurat package.  
E.g.

- Guided Clustering tutorial : We've just worked through this
- Differential expression : An Exploration of differential expression methods within Seurat
- Data integration : Seurat's data integration is a popular method to combine different datasets into one joint analysis.

### Seurat Cheatsheet

[https://satijalab.org/seurat/articles/essential\\_commands.html](https://satijalab.org/seurat/articles/essential_commands.html)

A useful resource for asking; How can I do 'X' with my seurat object?

### OSCA

<https://bioconductor.org/books/release/OSCA/>

An comprehensive resource for analysis approaches for single cell data. This uses the SingleCellExperiment bioconductor ecosystem, but alot of the same principle still apply.

This includes a good discussion of useing pseudobulk approaches, worth checking out for differential expression analyses.

### **MBP training Reading list**

<https://monashbioinformaticsplatform.github.io/Single-Cell-Workshop/>

A workshop page for a previous workshop (upon which this one is based) run by Monash Bioinformatics Platform - down the bottom there is an extensive list of useful single cell links and resources.

### **Biocommons Single Cell Omics**

<https://www.biocommons.org.au/single-cell-omics>

Join the single cell omics community resources being setup by biocommons.

---

## **15.2 Data**

### **Demo 10X data**

<https://www.10xgenomics.com/resources/datasets>

10X genomics have quite a few example datasets availble for download (including PBMC3k). This is a useful resource if you want to see what the ‘raw’ data looks like for a particular technology.

### **GEO**

<https://www.ncbi.nlm.nih.gov/geo/>

Many papers publish their raw single cell data in GEO. Formats vary, but often you can find the counts matrix. # (PART) Other resources {-}

### **Seurat data**

<https://github.com/satijalab/seurat-data>

Package for obtaining a few datasets as seurat objects.

---

### 15.3 Analysis Tools

A handful of the many tools that might be worth checking out for next steps.

#### Cyclone

<https://pubmed.ncbi.nlm.nih.gov/26142758/>

Part of the scran package, cyclone is a(nother) method for determining cell phase. Doco

#### Harmony

<https://portals.broadinstitute.org/harmony/articles/quickstart.html>

Method for integration of multiple single cell datasets.

#### SingleR

<http://bioconductor.org/books/release/SingleRBook/>

There is extensive documentation for the singleR package in the ‘singleR’ book.

#### Scrublet

<https://github.com/swolock/scrublet>

A python based tool for doublet detection. One of many tools in this space.

#### ScVelo

<https://scvelo.readthedocs.io/>

A package for single cell RNA velocity analysis, useful for developmental/pseudotime trajectories. Python/scanpy based.

#### Monocle

<https://cole-trapnell-lab.github.io/monocle3/>

A package for single cell developmental//pseudotime trajectory analysis.

#### TidySeurat

<https://stemangiola.github.io/tidyseurat/>

For fans of tidyverse-everything, there’s tidyseurat. Example workflow here

---

## 15.4 Preprocessing Tools

Tooks that process raw sequencing data into counts matricies

### Cell Ranger

<https://support.10xgenomics.com/single-cell-gene-expression/software/pipelines/latest/what-is-cell-ranger>

CellRanger is the 10X tool that takes raw fastq sequence files and produces the counts matricies that are the starting point for today's analysis. It only works for 10X data.

### STAR Solo

STAR is an aligner (which is actually used within cell ranger). STAR Solo is a tool for producing counts matricies, and is configurable enough for use with multiple technologies.

<https://github.com/alexdobin/STAR/blob/master/docs/STARsolo.md>



## **Part IV**

# **Seurat Object**



# Chapter 16

## Structure

### 16.1 Load an existing Seurat object

The data we're working with today is a small dataset of about 5000 PBMCs (peripheral blood mononuclear cells) from a healthy donor.

First, load Seurat package.

```
library(Seurat)
```

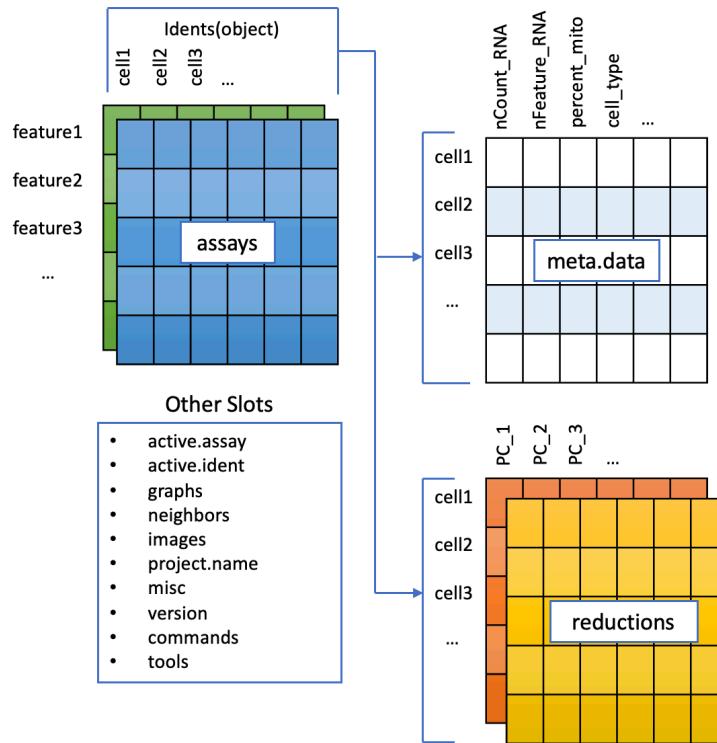
And here's the one we prepared earlier. Seurat objects are usually saved as 'rds' files, which is an R format for storing binary data (not-text or human-readable). The functions `readRDS()` can load it.

```
seurat_object <- readRDS("data/kang2018.rds")
```

#### Discussion: The Seurat Object in R

Lets take a look at the seurat object we have just created in R, `seurat_object`

To accomodate the complexity of data arising from a single cell RNA seq experiment, the seurat object keeps this as a container of multiple data tables that are linked.



The functions in seurat can access parts of the data object for analysis and visualization, we will cover this later on.

There are a couple of concepts to discuss here.

### Class

These are essentially data containers in R as a class, and can be accessed as a variable in the R environment.

Classes are pre-defined and can contain multiple data tables and metadata. For Seurat, there are three types.

- Seurat - the main data class, contains all the data.
- Assay - found within the Seurat object. Depending on the experiment a cell could have data on RNA, ATAC etc measured
- DimReduc - for PCA and UMAP

### Slots

Slots are parts within a class that contain specific data. These can be lists, data tables and vectors and can be accessed with conventional R methods.

### Data Access

Many of the functions in Seurat operate on the data class and slots within them

seamlessly. There maybe occasion to access these separately to `hack` them, however this is an advanced analysis method.

The ways to access the slots can be through methods for the class (functions) or with standard R accessor nomenclature.

### Examples of accessing a Seurat object

The `assays` slot in `seurat_object` can be accessed with `seurat_object@assays`.

The RNA assay can be accessed from this with `seurat_object@assays$RNA`.

We often want to access assays, so Seurat nicely gives us a shortcut `seurat_object$RNA`. You may sometimes see an alternative notation `seurat_object[["RNA"]]`.

In general, slots that are always in an object are accessed with `@` and things that may be different in different data sets are accessed with `$`.

### Have a go

Use `str` to look at the structure of the Seurat object `seurat_object`.

What is in the `meta.data` slot within your Seurat object currently? What type of data is contained here?

Where is our count data within the Seurat object?

## 16.2 What's in there?

Some of the most important information for working with Seurat objects is in the metadata. This is cell level information - each row is one cell, identified by its barcode. Extra information gets added to this table as analysis progresses.

```
head(seurat_object@meta.data)
#>                 orig.ident nCount_RNA nFeature_RNA ind
#> AGGGCGCTATTCC-1 SeuratProject      2020      523 1256
#> GGAGACGATTGTT-1 SeuratProject       840      381 1256
#> CACCGTTGTCGTAG-1 SeuratProject     3097      995 1016
#> TATCGTACACGCAT-1 SeuratProject     1011      540 1488
#> TGACGCCCTTGCTTT-1 SeuratProject     570      367 101
#> TACGAGACCTATTTC-1 SeuratProject    2399      770 1244
#>                      stim   cell multiplets
#> AGGGCGCTATTCC-1 stim   CD14+ Monocytes singlet
#> GGAGACGATTGTT-1 stim   CD4 T cells singlet
#> CACCGTTGTCGTAG-1 ctrl   FCGR3A+ Monocytes singlet
#> TATCGTACACGCAT-1 stim   B cells singlet
#> TGACGCCCTTGCTTT-1 ctrl   CD4 T cells ambs
#> TACGAGACCTATTTC-1 stim   CD4 T cells singlet
```

That doesn't have any gene expression though, that's stored in an 'Assay'. The Assay structure has some nuances (see discussion below), but there are functions that get the assay data out for you.

By default this object will return the normalised data (from the only assay in this object, called RNA). Every '?' is a zero.

```
GetAssayData(seurat_object)[1:15,1:2]
#> 15 x 2 sparse Matrix of class "dgCMatrix"
#>          AGGGCGCTATTCC-1 GGAGACGATTCGTT-1
#> MIR1302-10          .
#> FAM138A            .
#> OR4F5              .
#> RP11-34P13.7       .
#> RP11-34P13.8       .
#> AL627309.1        .
#> RP11-34P13.14      .
#> RP11-34P13.9       .
#> AP006222.2        .
#> RP4-669L17.10      .
#> OR4F29             .
#> RP4-669L17.2       .
#> RP5-857K21.15      .
#> RP5-857K21.1       .
#> RP5-857K21.2       .
```

But the raw counts data is accessible too.

```
GetAssayData(seurat_object, slot='counts')[1:15,1:2]
#> 15 x 2 sparse Matrix of class "dgCMatrix"
#>          AGGGCGCTATTCC-1 GGAGACGATTCGTT-1
#> MIR1302-10          .
#> FAM138A            .
#> OR4F5              .
#> RP11-34P13.7       .
#> RP11-34P13.8       .
#> AL627309.1        .
#> RP11-34P13.14      .
#> RP11-34P13.9       .
#> AP006222.2        .
#> RP4-669L17.10      .
#> OR4F29             .
#> RP4-669L17.2       .
#> RP5-857K21.15      .
#> RP5-857K21.1       .
#> RP5-857K21.2       .
```

## **Chapter 17**

# **Acknowledgements**

This material is mostly based on the Seurat introductory tutorial.

It also draws from material for a workshop provided by QCIF developed by Sarah Williams and Ahmed Mehdi for the Australian Biocommons.



# Chapter 18

## Session info

```
library(pander)
demo.Rmd_session <- sessionInfo()
pander(demo.Rmd_session)
```

R version 4.4.0 (2024-04-24)

Platform: x86\_64-pc-linux-gnu

locale: LC\_CTYPE=en\_AU.UTF-8, LC\_NUMERIC=C, LC\_TIME=en\_AU.UTF-8, LC\_COLLATE=en\_AU.UTF-8, LC\_MONETARY=en\_AU.UTF-8, LC\_MESSAGES=en\_AU.UTF-8, LC\_PAPER=en\_AU.UTF-8, LC\_NAME=C, LC\_ADDRESS=C, LC\_TELEPHONE=C, LC\_MEASUREMENT=en\_AU.UTF-8 and LC\_IDENTIFICATION=C

attached base packages: stats4, stats, graphics, grDevices, utils, datasets, methods and base

other attached packages: pander(v.0.6.6), edgeR(v.4.2.2), limma(v.3.60.6), celldex(v.1.14.0), SingleR(v.2.6.0), SingleCellExperiment(v.1.26.0), SummarizedExperiment(v.1.34.0), Biobase(v.2.64.0), GenomicRanges(v.1.56.2), GenomeInfoDb(v.1.40.1), IRanges(v.2.38.1), S4Vectors(v.0.42.1), BiocGenerics(v.0.50.0), MatrixGenerics(v.1.16.0), matrixStats(v.1.5.0), clustree(v.0.5.1), ggraph(v.2.2.1), harmony(v.1.2.3), Rcpp(v.1.1.0), future(v.1.67.0), patchwork(v.1.3.1), Seurat(v.5.3.0), SeuratObject(v.5.1.0), sp(v.2.2.0), ggplot2(v.3.5.2) and dplyr(v.1.1.4)

loaded via a namespace (and not attached): spatstat.sparse(v.3.1-0), httr(v.1.4.7), RColorBrewer(v.1.1-3), tools(v.4.4.0), settransform(v.0.4.2), backports(v.1.5.0), alabaster.base(v.1.4.2), utf8(v.1.2.6), R6(v.2.6.1), HDF5Array(v.1.32.1), lazyeval(v.0.2.2), uwot(v.0.2.3), rhdf5filters(v.1.16.0), withr(v.3.0.2), gridExtra(v.2.3), progressr(v.0.15.1), cli(v.3.6.5), spatstat.explore(v.3.5-2), fastDum-

*mies(v.1.7.5), alabaster.se(v.1.4.1), labeling(v.0.4.3), spatstat.data(v.3.1-6), ggridges(v.0.5.6), pbapply(v.1.7-4), parallelly(v.1.45.1), rstudioapi(v.0.17.1), RSQLite(v.2.4.2), generics(v.0.1.4), ica(v.1.0-3), spatstat.random(v.3.4-1), Matrix(v.1.7-0), ggbeeswarm(v.0.7.2), abind(v.1.4-8), lifecycle(v.1.0.4), yaml(v.2.3.10), rhdf5(v.2.48.0), SparseArray(v.1.4.8), BiocFileCache(v.2.12.0), Rtsne(v.0.17), grid(v.4.4.0), blob(v.1.2.4), promises(v.1.3.3), ExperimentHub(v.2.12.0), crayon(v.1.5.3), miniUI(v.0.1.2), lattice(v.0.22-6), beachmat(v.2.20.0), cowplot(v.1.2.0), KEGGREST(v.1.44.1), pillar(v.1.11.0), knitr(v.1.50), future.apply(v.1.20.0), codetools(v.0.2-20), glue(v.1.8.0), spatstat.univar(v.3.1-4), data.table(v.1.17.8), gypsum(v.1.0.1), vctrs(v.0.6.5), png(v.0.1-8), spam(v.2.11-1), gtable(v.0.3.6), cachem(v.1.1.0), xfun(v.0.52), S4Arrays(v.1.4.1), mime(v.0.13), tidygraph(v.1.3.1), survival(v.3.5-8), tinytex(v.0.57), statmod(v.1.5.0), fitdistrplus(v.1.2-4), ROCOCR(v.1.0-11), nlme(v.3.1-164), bit64(v.4.6.0-1), alabaster.ranges(v.1.4.2), filelock(v.1.0.3), RcppAnnoy(v.0.0.22), irlba(v.2.3.5.1), vapor(v.0.4.7), KernSmooth(v.2.23-22), colorspace(v.2.1-1), DBI(v.1.2.3), ggrastr(v.1.0.2), tidyselect(v.1.2.1), bit(v.4.6.0), compiler(v.4.4.0), curl(v.6.4.0), httr2(v.1.2.1), hdf5r(v.1.3.12), DelayedArray(v.0.30.1), plotly(v.4.11.0), bookdown(v.0.43), checkmate(v.2.3.2), scales(v.1.4.0), lmtest(v.0.9-40), rappdirs(v.0.3.3), stringr(v.1.5.1), digest(v.0.6.37), goftest(v.1.2-3), spatstat.utils(v.3.1-5), presto(v.1.0.0), alabaster.matrix(v.1.4.2), rmarkdown(v.2.29), XVector(v.0.44.0), RhpcBLASctl(v.0.23-42), htmltools(v.0.5.8.1), pkgconfig(v.2.0.3), sparseMatrixStats(v.1.16.0), dbplyr(v.2.5.0), fastmap(v.1.2.0), rlang(v.1.1.6), htmlwidgets(v.1.6.4), UCSC.utils(v.1.0.0), shiny(v.1.11.1), DelayedMatrixStats(v.1.26.0), farver(v.2.1.2), zoo(v.1.8-14), jsonlite(v.2.0.0), BiocParallel(v.1.38.0), BiocSingular(v.1.20.0), magrittr(v.2.0.3), GenomeInfoDbData(v.1.2.12), dotCall64(v.1.2), Rhdf5lib(v.1.26.0), viridis(v.0.6.5), reticulate(v.1.43.0), stringi(v.1.8.7), alabaster.schemas(v.1.4.0), zlibbioc(v.1.50.0), MASS(v.7.3-60.2), AnnotationHub(v.3.12.0), plyr(v.1.8.9), parallel(v.4.4.0), listenv(v.0.9.1), ggrepel(v.0.9.6), deldir(v.2.0-4), Biostrings(v.2.72.1), graphlayouts(v.1.2.2), splines(v.4.4.0), tensor(v.1.5.1), locfit(v.1.5-9.12), igraph(v.2.1.4), spatstat.geom(v.3.5-0), RcppHNSW(v.0.6.0), reshape2(v.1.4.4), ScaledMatrix(v.1.12.0), BiocVersion(v.3.19.1), evaluate(v.1.0.4), BiocManager(v.1.30.26), tweenr(v.2.0.3), httpuv(v.1.6.16), RANN(v.2.6.2), tidyverse(v.1.3.1), purrr(v.1.1.0), polyclip(v.1.10-7), scattermore(v.1.2), ggforce(v.0.5.0), rsrvd(v.1.0.5), xtable(v.1.8-4), RSpectra(v.0.16-2), later(v.1.4.2), viridisLite(v.0.4.2), tibble(v.3.3.0), memoise(v.2.0.1), beeswarm(v.0.4.0), AnnotationDbi(v.1.66.0), cluster(v.2.1.6) and globals(v.0.18.0)*