

Core Steps in scRNA-seq Analysis: QC to Clustering

Single-Cell & Spatial Omics Series · NCI CCR BTEP · 2026

Alex Emmons, PhD


2026-05-20

Single-Cell & Spatial Omics 2026

Upcoming Events

- **Core Steps in scRNA-seq Analysis: QC to Clustering** ← *you are here*
- Multi-sample analysis in single cell RNASeq: Batch correction, annotation, and differential expression (5/27)
- Multimodal Spatial Transcriptomic Analysis (6/03)
- From One Cell to Multiple Molecular Views: A Practical Introduction to Single-Cell Multi-omics (6/10)

Resources

-  [Series Home](#)
-  [Download Tutorial Data](#)
-  ncibtep@nih.gov

Learning Objectives

By the end of this session you will be able to:

1. **Apply quality control filters** to retain high-quality cells
2. **Normalize** data using `SCTransform()`
3. **Perform dimensionality reduction** with PCA
4. **Cluster** cells and visualize with UMAP

Warning

Thresholds and clustering decisions shown here are **starting points** – revisit them in the context of your biology.

The Dataset

Lung regeneration after influenza infection [Niethamer et al. 2025, Cell Stem Cell](#)

Sample	Condition	Time Point
S89Inf42	Influenza	Day 42 post-infection
S91Hom42	Homeostasis	Day 42
S93Hom3	Homeostasis	Day 3
S99Inf42	Influenza	Day 42 post-infection

- **2 biological replicates** of the infected condition at day 42
- Comparison of injury-associated states vs. homeostatic controls
- Expect both **technical** and **biological** variation between samples

Load Packages & Object

```
1 library(Seurat)
2 library(tidyverse)
3 library(hdf5r)
4 library(patchwork)
5 library(scuttle)
6 library(glmGamPoi) # speeds up SCTransform
7 library( presto)   # speeds up FindAllMarkers
```

```
1 lung <- readRDS("./outputs/merged_Seurat_lung.rds")
2 lung
3 Layers(lung)
```

Note

In **Seurat v5**, merged RNA assays retain per-sample matrices as separate layers (e.g., `counts.S89Inf42`). This enables per-sample QC even on a merged object.

Quality Control: The Goal

Remove **low-quality cells** before downstream analysis

`nCount_RNA`

Total UMI count per cell

- ▲ High → possible doublet
- ▼ Low → ambient RNA / empty droplet

`nFeature_RNA`

Genes detected per cell

- ▲ High → possible doublet
- ▼ Low → ambient RNA / low complexity

`percent_mt`

Mitochondrial read fraction

- ▲ High → dying/damaged cell

⚠ Always consider **biological explanations** — high mt% can reflect metabolic activity; low counts may indicate quiescent populations.

Add Mitochondrial Percentage

```
1 # Mouse dataset – mitochondrial genes start with 'mt-'  
2 # Human datasets use '^MT-'  
3 lung[["percent.mt"]] <- PercentageFeatureSet(lung, pattern = "^mt-")
```

```
1 # Extract metadata for plotting  
2 metadata <- lung[[[]]]  
3 head(metadata)
```


Note

QC should be sample-aware. Visualize metrics per sample – differences may reflect biology, not just technical issues.


QC Assessment Strategy

Metrics should be assessed jointly, not in isolation

Manual thresholds

- Subjective visual assessment
- Can be arbitrary
-  What we do here

Adaptive thresholds

- Uses the data's own distribution
- Outliers defined as 3 or 5 MAD from the median
- `scuttle::isOutlier()`
-  OSCA: QC outlier detection

Our Strategy:

- **Step 1: Visualize QC metrics**
- **Step 2: Decide on thresholds**
- **Step 3: Filter the Object**

Step 1: Visualize QC Metrics

```
1 VlnPlot(lung, features = "nCount_RNA", layer = "counts",
2         group.by = "orig.ident", alpha = 0.05)
3
4 metadata %>%
5   ggplot(aes(x = nCount_RNA, fill = orig.ident, color = orig.ident)) +
6   geom_density(alpha = 0.2) + theme_classic() + scale_x_log10() +
7   geom_vline(xintercept = 600, color = "red", linetype = "dotted")
```

```
1 metadata %>%
2   ggplot(aes(x = nFeature_RNA, fill = orig.ident, color = orig.ident)) +
3   geom_density(alpha = 0.2) + theme_classic() + scale_x_log10() +
4   geom_vline(xintercept = c(375, 5500), color = "red", linetype = "dotted")
5
6 metadata %>%
7   ggplot(aes(x = percent.mt, fill = orig.ident, color = orig.ident)) +
8   geom_density(alpha = 0.2) + theme_classic() + scale_x_log10() +
9   geom_vline(xintercept = 15)
```

Check out [scCustomize](#) for additional plotting options and flexible color palettes.

Step 2 / 3: Set QC Thresholds and Filter

Chosen thresholds (conservative starting point):

Metric	Filter
nFeature_RNA	> 375
nCount_RNA	> 600
percent.mt	< 15%

```
1 keep_cells <- metadata %>%
2   rownames_to_column("cell") %>%
3   filter(nFeature_RNA > 375,
4         nCount_RNA > 600,
5         percent.mt < 15) %>%
6   pull(cell)
7
8 lung_filt <- subset(lung, cells = keep_cells)
```

Check: Cells Before & After Filtering

```
1 bind_rows(  
2   before = metadata,  
3   after = lung_filt[[]],  
4   .id = "stage"  
5 ) %>%  
6   count(stage, orig.ident) %>%  
7   ggplot(aes(x = orig.ident, y = n, fill = stage)) +  
8   geom_col(position = "dodge") + theme_classic() +  
9   geom_text(aes(label = n, group = stage),  
10            position = position_dodge(width = 0.9))
```

```
1 saveRDS(lung_filt, "./outputs/merged_Seurat_lung_qcfiltered.rds")
```



Tip



Save your filtered object – it's a useful checkpoint to return to.

Why Normalize?

Raw counts are **not directly comparable** across cells

- Cells vary in total mRNA content and technical capture efficiency
- Sources of technical variation include differences in cell lysis, reverse transcription efficiency, and stochastic molecular sampling
- Normalization minimizes technical effects so expression can be compared across cells

Standard workflow

1. `NormalizeData()`
2. `FindVariableFeatures()`
3. `ScaleData()`

Modern approach

- ✓ `SCTransform()` — one step, better statistical model

SCTransform

```
1 lung_filt <- SCTransform(  
2   lung_filt,  
3   vars.to.regress = "percent.mt",  
4   verbose = FALSE  
5 )
```

SCTransform output — **SCT** assay layers:

Layer	Contents
<code>counts</code>	Regularized corrected counts
<code>data</code>	Log-normalized corrected counts
<code>scale.data</code>	Pearson residuals for variable features → used in PCA

```
1 DefaultAssay(lung_filt) # should now be "SCT"
```

Dimensionality Reduction: PCA

Single-cell data is **highly dimensional** — thousands of genes per cell

PCA summarizes correlated gene expression patterns into principal components (PCs), which are used as the basis for clustering

```
1 lung_filt <- RunPCA(lung_filt, assay = "SCT", verbose = FALSE)
```

```
1 ElbowPlot(lung_filt, ndims = 40)
```

```
1 DimHeatmap(lung_filt, dims = 1:9, cells = 500, balanced = TRUE, ncol = 3)
```

```
2 DimHeatmap(lung_filt, dims = 20:30, cells = 500, balanced = TRUE, ncol = 3)
```

Graph-Based Clustering

Three steps:

1. **FindNeighbors()** – builds a KNN graph in PCA space; refines edge weights using Jaccard similarity
2. **FindClusters()** – partitions the graph into communities using modularity optimization
3. **RunUMAP()** – 2D visualization of the structure

Resolution parameter

Controls granularity

▲ Higher → more clusters

▼ Lower → fewer clusters

Try several and compare!

Run Clustering & UMAP

```
1 lung_filt <- FindNeighbors(lung_filt, dims = 1:30)
2 lung_filt <- FindClusters(lung_filt, resolution = 0.2)
3 lung_filt <- RunUMAP(lung_filt, dims = 1:30)
```

```
1 DimPlot(
2   lung_filt, reduction = "umap",
3   group.by = c("orig.ident", "seurat_clusters",
4               "time_point", "treatment"),
5   ncol = 2, alpha = 0.2
6 )
7
8 DimPlot(lung_filt, reduction = "umap",
9         group.by = "orig.ident",
10        split.by = "treatment", ncol = 2, alpha = 0.2)
```

Annotate with Canonical Markers

```
1 FeaturePlot(lung_filt, features = c("Epcam", # epithelial
2                                     "Ptprc")) # immune (CD45)
```

Canonical markers are biologically validated genes known to define specific cell types

Combine with `VlnPlot()`, `DotPlot()` for detailed inspection

Annotation tools:

- Azimuth
- SingleR
- scType
- Garnett

Find Cluster Markers

```
1 # Required before DE testing with SCT + multiple samples
2 lung_filt <- PrepSCTFindMarkers(lung_filt, verbose = TRUE)







1 lung_filt_markers <- FindAllMarkers(lung_filt, only.pos = TRUE)
2
3 lung_filt_markers <- lung_filt_markers |>
4   group_by(cluster) |>
5   arrange(desc(avg_log2FC), pct.1, p_val_adj, .by_group = TRUE)
```

Column	Meaning
pct.1	% of cells in this cluster expressing the gene
pct.2	% in all other clusters
avg_log2FC	Fold change vs. other clusters

Save & Next Steps

```
1 saveRDS(lung_filt, "./outputs/merged_Seurat_lung_sct_clustered.rds")
```

Typical next steps:

-  **Integration** — if cells cluster primarily by sample or batch
-  **Cluster annotation** — markers + reference tools + biology
-  **Subsetting & re-clustering** — zoom into populations of interest
-  **Differential expression** — across clusters, conditions, or time points
-  **Cell composition analysis** — do proportions differ across treatments?
-  **Trajectory analysis** — model injury → regeneration state transitions

Resources & References

Tutorial & Data

- [Series home](#)
- [Seurat v5 tutorial](#)
- [SCTransform vignette](#)
- [Single-Cell Best Practices](#)
- [HBC scRNA-seq training](#)

Data

- [Niethamer et al. 2025, *Cell Stem Cell*](#)

Key Papers

- Luecken & Theis (2019) *Mol Syst Biol* — best practices review
- Heumos et al. (2023) *Nat Rev Genet* — best practices across modalities
- Hafemeister & Satija (2019) *Genome Biol* — SCTransform

Questions?

 ncibtep@nih.gov

 Office hours & video archive at bioinformatics.ccr.cancer.gov

Appendix: Adaptive QC with MAD

An alternative to manual thresholds – `scuttle::isOutlier()`

```
1 library(scuttle)
2
3 outlier_low_feat <- isOutlier(
4   metadata$nFeature_RNA,
5   nmads = 3,
6   type = "lower",
7   log = TRUE
8 )
9 summary(outlier_low_feat)
```

Outliers are generally defined as 3 or 5 MADs from the median. More information on adaptive thresholds: [OSCA: QC outlier detection](#)

Appendix: Standard Normalization Workflow

For comparison — the classic three-step approach:

```
1 lung_filt <- NormalizeData(lung_filt,  
2     normalization.method = "LogNormalize",  
3     scale.factor = 10000)  
4  
5 lung_filt <- FindVariableFeatures(lung_filt,  
6     selection.method = "vst",  
7     nfeatures = 2000)  
8  
9 lung_filt <- ScaleData(lung_filt,  
10     vars.to.regress = "percent.mt")
```

We prefer `SCTransform()` for this dataset, but both are valid starting points.