Analysing Single-Cell RNA-Seq with R

v2025-06 (Seurat v5)

Simon Andrews simon.andrews@babraham.ac.uk



Major scRNA Package Systems



https://satijalab.org/seurat/

scater

Single-Cell Analysis Toolkit for Gene Expression Data in R

https://bioconductor.org/packages/release/bioc/html/scater.html



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





https://scanpy.readthedocs.io/en/stable/

What do they provide?

- scRNA Data Structure
 - Data parsing
 - Quantitations
 - Metadata
- Processing methods
 - Normalisation
 - Integration
 - Dimensionality reduction
 - PCA / tSNE / UMAP

- Statistics
 - Enriched genes
 - Differential expression

- Plotting
 - Projections
 - QC
 - Expression graphs

Seurat

- Most popular R framework
 - Well supported and frequently updated
- Good documentation
 - Reference documentation
 - Vignettes and examples
- Easy data model to work with
- Lots of built in functionality
 - Easy to extend to build your own

Seurat 5.2.0



Seurat v5

We are excited to release Seurat v5! To install, please follow the instructions in our install page. This update brings the following new features and functionality:

• Integrative multimodal analysis: The cellular transcriptome is just one aspect of cellular identity, and recent technologies enable routine profiling of chromatin accessibility, histone modifications, and protein levels from single cells. In Seurat v5, we introduce 'bridge integration', a statistical method to integrate experiments measuring different modalities (i.e. separate scRNA-seq and scATAC-seq datasets), using a separate multiomic dataset as a molecular 'bridge'. For example, we demonstrate how to map scATAC-seq datasets onto scRNA-seq datasets, to assist users in interpreting and annotating data from new modalities.

We recognize that while the goal of matching shared cell types across datasets may be important for many problems, users may also be concerned about which method to use, or that integration could result in a loss of biological resolution. In Seurat v5, we also introduce flexible and streamlined workflows for the integration of multiple scRNA-seq datasets. This makes it easier to explore the results of different integration methods, and to compare these results to a workflow that excludes integration steps.

- · Paper: Dictionary learning for integrative, multimodal, and scalable single-cell analysis
- Vignette: Streamlined integration of scRNA-seq data
- Vignette: Cross-modality bridge integration
- Website: Azimuth-ATAC, reference-mapping for scATAC-seq datasets

Seurat Data Structure

- Single object holds all data
 - Build from text table or 10X output (feature matrix h5 or raw matrix)



tools

list [0]

Seurat Metadata

Cell.Barcode <chr></chr>	orig.ident <fctr></fctr>	nCount_RNA <dbl></dbl>	nFeature_RNA <int></int>	percent.MT <dbl></dbl>	percent.Ribosomal <dbl></dbl>
AAACCTGAGAAACCAT-1	course	4410	1176	3.1519274	47.346939
AAACCTGAGATAGCAT-1	course	4470	1482	4.8098434	27.494407
AAACCTGAGCGTGAAC-1	course	2298	870	3.1331593	35.161010
AAACCTGCAACGCACC-1	course	4307	1433	4.2953332	26.259577
AAACCTGCACACCGAC-1	course	2001	669	2.6486757	42.528736
AAACCTGCATCGATTG-1	course	3310	918	3.7764350	48.791541
AAACCTGTCCAGATCA-1	course	3005	1032	2.9284526	38.735441
AAACGGGAGCGCTTAT-1	course	6932	1906	5.5972302	35.776111
AAACGGGAGCTACCTA-1	course	3714	1363	4.2003231	22.240172
AAACGGGAGTGATCGG-1	course	8020	2113	4.3266833	34.501247

• QC

Conditions

Clusters

data[[]]
data\$nCount_RNA
new data -> data\$new_metric

Access whole table Access one column Add a column

Seurat Quantitative Data

> LayerData(data, layer="counts")

Gene <chr></chr>	AACTCAGAGATGTAAC-1 <dbl></dbl>	AACTCAGAGTCCAGGA-1 <dbl></dbl>	AACTCAGCAAGTAATG-1
AL627309.1	0	0	0
AL627309.5	0	0	0
LINC01409	0	0	0
LINC01128	0	1	0
LINC00115	0	0	0
FAM41C	0	0	0
NOC2L	1	0	0

> LayerData(data, layer="data")

Gene <chr></chr>	AACTCTTAGCCGGTAA-1 <dbl></dbl>	AACTCTTCACTGTTAG-1 <dbl></dbl>	AACTCTTCATGCCTAA-1
AL627309.1	0.000000	0.000000	0
AL627309.5	0.000000	0.000000	0
LINC01409	0.000000	0.000000	0
LINC01128	0.000000	0.000000	0
LINC00115	0.000000	0.000000	0
FAM41C	0.000000	0.000000	0
NOC2L	1.228141	1.567888	0

Seurat Dimensionality Reductions

> Embeddings(data,reduction = "pca")

	PC_1	PC_2	PC_3
GAGAAACCAT-1	-4.135099	-8.7585629	0.004920869
GAGATAGCAT-1	10.279728	0.8451562	0.287589575
GAGCGTGAAC-1	-6.002598	4.9504875	-3.022266598
GCAACGCACC-1	10.610838	0.4030928	0.165408128
GCATCGATTG-1	-5.581052	-0.4542359	5.166186308

> Loadings (data, reduction="pca")

	PC_1	PC_2	PC_3
IGLC3	-0.01326925	-0.07255036	-0.017863727
IGLC1	-0.01039493	-0.05451957	-0.008700478
IGLC2	-0.01569606	-0.08276587	-0.021715297
IGKC	-0.02226017	-0.12668741	-0.038651318
S100A9	0.11325462	0.01836318	0.013403732

Variable Gene Information

> HVFInfo(data)

Gene <chr></chr>	mean <dbl></dbl>	variance <dbl></dbl>	variance.expected <dbl></dbl>	variance.standardized <dbl></dbl>
AL627309.1	5.077431e-04	5.076141e-04	5.099960e-04	0.9953296
AL627309.5	2.792587e-03	2.785496e-03	2.907594e-03	0.9580070
LINC01409	2.564103e-02	2.498991e-02	2.679890e-02	0.9324974
LINC01128	1.624778e-02	1.801934e-02	1.712920e-02	1.0519658
LINC00115	5.077431e-03	5.560805e-03	5.344264e-03	1.0405185
FAM41C	9.647119e-03	9.556478e-03	1.016386e-02	0.9402406
NOC2L	1.381061e-01	1.342992e-01	1.519011e-01	0.8841229
KLHL17	2.792587e-03	2.785496e-03	2.907594e-03	0.9580070
PLEKHN1	4.061945e-03	4.554345e-03	4.262279e-03	1.0685233
HES4	1.345519e-02	1.987986e-02	1.419635e-02	1.4003494

Seurat Methods

- Data Parsing
 - Read10X
 - Read10X_h5*
 - CreateSeuratObject
- Data Normalisation
 - NormalizeData
 - ScaleData
- Graphics
 - Violin Plot metadata or expression (VlnPlot)
 - Feature plot (FeatureScatter)
 - Projection Plot (DimPlot, DimHeatmap)

- Dimension reduction
 - RunPCA
 - RunTSNE
 - RunUMAP
- Statistics
 - Select Variable Genes
 FindVariableFeatures
 - Build nearest neighbour graph FindNeighbors
 - Build graph based cell clusters FindClusters
 - Find genes to classify clusters (multiple tests)
 FindMarkers

*Requires installing the hdf5r package



Example 10X Seurat Workflow



Example Seurat Workflow



Reading Data

Read10x_h5("filtered_feature_bc_matrix.h5") -> data

CreateSeuratObject(counts=data, project="course",) -> data

> data

An object of class Seurat 17136 features across 3939 samples within 1 assay Active assay: RNA (17136 features, 500 variable features)

- 3 layers present: counts, data, scale.data
- 2 dimensional reductions calculated: pca, tsne

Reading Multiple Files

CreateSeuratObject(Read10X_h5("SAMPLE1.h5"),project="sample1") -> sample1 CreateSeuratObject(Read10X_h5("SAMPLE2.h5"),project="sample2") -> sample2 CreateSeuratObject(Read10X_h5("SAMPLE3.h5"),project="sample3") -> sample3

```
merge(
```

```
sample1,
c(sample2,sample3),
add.cell.ids=c("sample1","sample2","sample3")
) -> data
```

	orig.ident <chr></chr>	nCount_RNA <dbl></dbl>	nFeature_RNA <int></int>	
sample1_AAACCTGAGCTGTTCA-1	sample 1	2216	1089	
sample1_AAACCTGCAACACGCC-1	sample 1	9103	2881	
sample1_AAACCTGCATGCAACT-1	sample 1	3927	1459	
sample1_AAACCTGGTGGTCTCG-1	sample1	6434	2251	
sample1_AAACCTGTCATCGATG-1	sample 1	2623	1318	
sample1_AAACGGGCATCCTTGC-1	sample 1	4903	1948	
sample1_AAACGGGGTAGAGTGC-1	sample 1	2727	934	
sample1_AAACGGGGTTGAACTC-1	An object of class s	Seurat		
sample1_AAAGATGTCACAGGCC-1	33696 features acros	33696 features across 6302 samples within 1 assav		
sample1_AAAGCAAAGCGATCCC-1	Active assav: RNA (33696 features. Ø varia	ble features)	
	3 layers present:	counts.sample1, counts.	sample2, counts.s	

QC – What problems are likely?

- Lysed cells
- Dead or dying cells
- Empty GEMs
- Double (or more) occupied GEMs
- Cells in different cell cycle stages



Lysed Cells

- Outer membrane is ruptured cytoplasmic RNAs leak out
 - Loss of mature RNA, increase in pre-mRNA
 - Lower overall counts/features
 - Increase in nuclear RNAs
 - MALAT1 is an easy marker to use



Dead or Dying Cells

Cells undergoing apoptosis have very different transcriptomes

-Lower total RNA production

-Huge upregulation of mitochondrial transcription



Empty GEMs

- GEMs containing no cell will still produce some sequence
 - Background RNA in the flow medium
 - Will be worse with higher numbers of lysed cells

• Total amount of signal will be greatly reduced

• Will often cluster together



Double occupied GEMs

- Will get a mixed signal from two different cells
- Not as obvious a signal as empty GEMs
 - More UMIs/Features per cell
 - Intermediate clustering



QC and Cell Filtering

- Standard QC Measures
 - Number of observed genes per cell
 - Number of reads per cell
 - Relationship between the two
- Calculated QC Measures
 - Amount of mitochondrial reads
 - Amount of ribosomal reads
 - Marker genes (eg MALAT1)
 - Cell cycle



Adding QC Metrics

PercentageFeatureSet(data, pattern="^MT-") -> data\$percent.MT

- Mitochondrial
- Immunoglobulin
- Malat1
- Ribosomal

QC and Cell Filtering



nCount_RNA

Applying Filters

subset(
 data,
 nFeature_RNA>750 &
 nFeature_RNA < 2000 &
 percent.MT < 10
) -> data

Count Normalisation and Scaling

- Raw counts are biased by total reads per cell
- Counts are more stable on a log scale
- Standard normalisation is just log reads per 10,000 reads
- For PCA counts scale each gene's expression to a z-score
 Can also use this step to try to regress out unwanted effects

Count Normalisation and Scaling

```
NormalizeData(
    data,
    normalization.method = "lognormalize"
) -> data
```

ScaleData(data) -> data

Variable Feature Selection

- Selects a subset of genes to use for downstream analysis
- Identify genes with an unusual amount of variability
- Link the variability with the expression level to find variation which is high in the context of the expression level
- Keep only the most variable genes

FindVariableFeatures(
 data,
 selection.method = "vst",
 nfeatures=500
) -> data



Variable Features for Multiple Samples

• Can just find variable features across all cells

```
lapply(unique(data$orig.ident), function(x) {
    data |>
    subset(orig.ident==x) |>
    FindVariableFeatures() |>
    HVFInfo() |>
    as_tibble(rownames="Gene") |>
    add column(orig.ident=x)}) -> variability data
```

do.call(bind_rows,variability_data) -> variability_data

- Compare Variance
- Select Common Variable Features

Dimensionality Reduction

- Start with PCA on the normalised, filtered (both cells and genes), scaled data
- Scree / Elbow plot to decide how many PCs are informative
- Pass only the interesting PCs to subsequent tSNE or UMAP reduction to get down to 2 dimensions



Dimensionality Reduction





Defining clusters

- Construct nearest neighbour graph
 - Constructed from PCA
 - Same dimensions as tSNE/UMAP
- Find clusters
 - All cells are classified
 - Graph Based (Louvain) Clustering
 - Resolution (0.01 5) defines granularity



data, resolution = 0.5) -> data









Clustree to see effect of resolution



12 17 11

https://github.com/lazappi/clustree

Comparing Properties of Clusters

VlnPlot(data,features="nFeature RNA")

- We want to know that clusters are occurring because of biological changes, not technical differences
- We plot QC metrics for clusters
 - Read/Gene counts
 - Mitochondrion
 - MALAT1
- Can remove suspect clusters



subset(data, !Seurat_clusters %in% c(8,10,12)) -> data

Statistical Analysis



- Cluster 1 vs Clusters [2,3,4]
- Cluster 1 vs Cluster 3

• Cluster A1 vs Cluster B1

Types of Statistics

- Non Parametric Stats
 - Default test
 - Wilcoxon Rank Sum
 - Semi Quantitative
 - Each Cell is a replicate
 - Highly Powered
 - Over powered?

- ROC Analysis
 - Not a stats test
 - How well can each gene separate the groups
 - Values 0 1
 - 0.5 is worst
 - 0 or 1 are perfect separation

DESeq Stats

- Only with replicates
- Aggregate counts per cluster and sample
- Standard RNA-Seq analysis
- Works well if your data supports it

Statistical analysis of differences between clusters

```
FindMarkers(
    data,
    ident.1 = 2,
    ident.2 = 6,
    test.use = "roc",
    only.pos = TRUE
```

FindAllMarkers(
 data,
 group.by ="seurat_clusters",
 test.use = "roc",
 only.pos = TRUE





BMC Bioinformatics

RESEARCH ARTICLE

Open Access

CrossMark



Tianyu Wang¹, Boyang Li², Craig E. Nelson³ and Sheida Nabavi^{4*}

Conclusions: In general, agreement among the tools in calling DE genes is not high. There is a trade-off between true-positive rates and the precision of calling DE genes. Methods with higher true positive rates tend to show low precision due to their introducing false positives, whereas methods with high precision show low true positive rates due to identifying few DE genes. We observed that current methods designed for scRNAseq data do not tend to show better performance compared to methods designed for bulk RNAseq

Automated Cell Assignment

- Can automatically assign cell identities to clusters
- Need a source of marker genes
 - Result of a previous run/experiment
 - Publicly available data (https://azimuth.hubmapconsortium.org/)

- Many packages to do this
 - SCINA has worked well for us
 - Azimuth built into Seurat

Abdelaal *et al. Genome Biology* (2019) 20:194 https://doi.org/10.1186/s13059-019-1795-z

Genome Biology

RESEARCH

A comparison of automatic cell identification methods for single-cell RNA sequencing data

Tamim Abdelaal^{1,2†}, Lieke Michielsen^{1,2†}, Davy Cats³, Dylan Hoogduin³, Hailiang Mei³, Marcel J. T. Reinders^{1,2} and Ahmed Mahfouz^{1,2*}



Open Access

Integrating Multiple Runs

 When multiple runs are combined (eg Unstim and Stim), the batch differences between the runs can overwhelm the biological differences

 Raw comparisons can therefore miss changes between what are actually matched subgroups

Raw merged runs

• Two PBMC populations run at different times

 tSNE spread coloured by library

• Little to no overlap between cell populations



Integrating Runs

• Split the layers based on the metadata

split(data[["RNA"]], f = data\$Batch) -> data[["RNA"]]

- Rerun Normalisation, Variable Features, Scaling, PCA
- Create a new integrated layer

```
IntegrateLayers(
   object = data, method = RPCAIntegration,
   orig.reduction = "pca", new.reduction ="integrated.rpca",
   verbose = FALSE
) -> data
```

Integrating Runs



Over-Integration



Exercise – Using Seurat to analyse 10X data