

Single-cell RNA-seq data analysis in Chipster, 4.-5.3.2025

chipster@csc.fi

PART I: One sample analysis - Finding clusters of cells and marker genes for them

In this tutorial we detect subgroups of peripheral blood mononuclear cells (PBMCs), and we also want to find marker genes for the different cell types. The 10X Genomics data set used in the exercises is available at https://satijalab.org/seurat/articles/pbmc3k_tutorial.html. The tar package containing the three 10X Genomics output files has been already imported in Chipster for you.

Open Chipster: Go to <https://chipster.csc.fi/>, click on **Launch Chipster**, and log in.

1. Open training session

Click **Sessions**, go to **Training sessions** and select **course_single_cell_RNAseq_Seurat**. Rename the session **course_single_cell_RNAseq_Seurat_your_first_name**

2. Setup Seurat object & perform quality control

Select the **files.tar.gz** and the tool **Single-cell RNA-seq / Seurat v5 -Setup and QC**. Check the parameters, and set **Project name for plotting = PBMC**. Run the tool. Select the **QCplots.pdf** and click **Open in new tab**. Look at all the pages.

- What would be the optimal limits for the number of genes (nFeature_RNA) and mitochondrial transcript percentage (percent.mt)?
- How many cells are there?

3. Filter cells

Select **setup_seurat_obj.Robj** and the tool **Seurat v5 – Filter cells**. Are the default cell filtering parameters good for this dataset, based on the QC plots?

- How many cells were filtered out?

4. Normalize expression values, scale data, regress out unwanted variation, and detect highly variable genes

Select **seurat_obj_filter.Robj** and the tool **Seurat v5 – Normalize, regress and detect variable genes**. While the tool is running, click **Info** to open the tool manual page and learn about the steps this tool performs.

- What are those steps?

Once the tool is done, select the file **Dispersion_plot.pdf** and click **Open in new tab**. Check also the second page.

- What are the ten most highly variable genes?

5. Principal component analysis

Select **seurat_obj_preprocess.Robj** from the previous step and run the tool **Seurat v5-PCA** so that you set **Number of PCs to compute = 20**.

Select **PCAplots.pdf** and click **Open in new tab**. Look at the PC heatmaps and the elbow plot.

- How many principal components should we use for clustering? Would 10 be ok?

6. Clustering

Select **seurat_obj_pca.Robj** from the previous step and run the tool **Seurat v5 -Clustering** using the following parameters:

Number of principal components to use = 10

Resolution for granularity = 0.5

-Open **clusterPlot.pdf**. Does the coloring (= clustering) match the grouping found by tSNE and UMAP? How many clusters are there?

7. Detection of cluster marker genes

Select **seurat_obj_clustering.Robj** from the previous step and the tool **Seurat v5 -Find differentially expressed genes between the clusters**. In the parameters, set the parameters as indicated below and run the tool.

Find all markers = FALSE

Cluster of interest = 3

Limit testing to genes which are expressed in at least this fraction of cells = 0.25

Check which markers show higher than 4-fold difference in expression between cluster 3 and all other cells. Select **markers.tsv** and run the tool **Filter table by column value** from the **Utilities category** using the following parameters:

Column to filter by = avg_log2FC

Does the first column lack a title = yes

Cutoff = 2 (why do we put 2 here if we want a 4-fold difference?)

Filtering criteria = larger-than

-How many genes do you get?

8. Visualize markers

Choose **seurat_obj_clustering.Robj** generated in step 6. Select tool **Seurat v5 -Visualize genes**. Type a marker **gene name(s)** in the parameter field. Try for example with MS4A1, LYZ and PF4. You can enter several gene names at the same time, separated by comma (.). Set the parameters

Add labels on top of clusters in plot = yes

Plotting order of cells based on expression = yes

For each gene, list the average expression and percentage of cells expressing it in each cluster = yes

-Are the genes you selected good markers and for which clusters (check both the plots and the tables)?

9. Annotate clusters

Choose **seurat_obj_clustering.Robj** generated in step 6. Run tool **SingleR cluster annotation**.

-Open **SingleR_annotations_plots.pdf** and see how the clusters are annotated. What are the cells in cluster 3?

10. Rename clusters

Select **seurat_obj_clustering.Robj** generated in step 6 and **cluster_names.tsv** which contains cluster annotations based on marker genes from the Seurat vignette. Check that the files are correctly assigned. Run the tool **Seurat v5 -Rename clusters** and open the result

file **clusterPlotRenamed.pdf**. How well do these annotations match with what you got from SingleR in the previous exercise?

11. Color named clusters based on mitochondrial transcript percentage

Choose **seurat_obj_renamed.Robj** generated in step 10. Select tool **Seurat v5 -Visualise features in UMAP plot** and set

Feature = percent.mt

Add labels on top of clusters in plot = yes

-Are the clusters evenly colored? Is this what you would expect?

12. Subset based on gene expression

Choose again **seurat_obj_clustering.Robj** generated in step 6. Run tool **Seurat v5- Subset Seurat objects based on gene expression** for gene MS4A1.

Then select **seurat_obj_subset.Robj** and run the tool **Seurat v5 -Extract information from Seurat object**. Open the result file **slots.txt**.

-How many cells are left in the subset?

13. Share a session with a colleague (in this case with Eija and Maria)

Make sure that no file is selected. Go to the **Session info** panel, click the **three dots** next to the session name, and select **Share**. In the new window that opens

-click **Add rule**.

-In the **UserID** field, enter **jaas/demo**

-set **Rights = Read-only** (you don't want us to mess up your session!)

-Click **Save**.

-Click **Close**.

Check what is your own UserID: Click on your **username** (top right corner) and select **Account**.

14. Bonus exercise: Repeat the analysis with SCTransform

Repeat steps 4-6, but this time, use the tool **Seurat v5 -SCTransform: Normalize, regress and detect variable genes** for normalization. Make the following changes:

In step 5, set **Number of PCs to compute = 50**.

In step 6, set

Normalisation method used previously = SCTransform

Number of PCs to use = 30

Resolution = 0.8

-How many clusters are found now? Which cluster seems to correspond to the cluster 3 obtained with the global scaling normalization previously?

Repeat step 7 but look for marker genes for cluster **2**.

-How many marker genes do you get? Are they the same genes as what you got for cluster 3 earlier? (Tip: use the Venn diagram to compare the markers to those found when the global scaling normalization was used)?