

Final Project OMICS

Chairetaki Antonia

December 2025

Seurat Guided Clustering Tutorial

Figure 1: The following plot displaying the distribution of key quality control metrics per cell (Peripheral Blood Monoluclear cell), specifically the number of unique genes detected, total RNA molecule counts, and the percentage of mitochondrial reads.

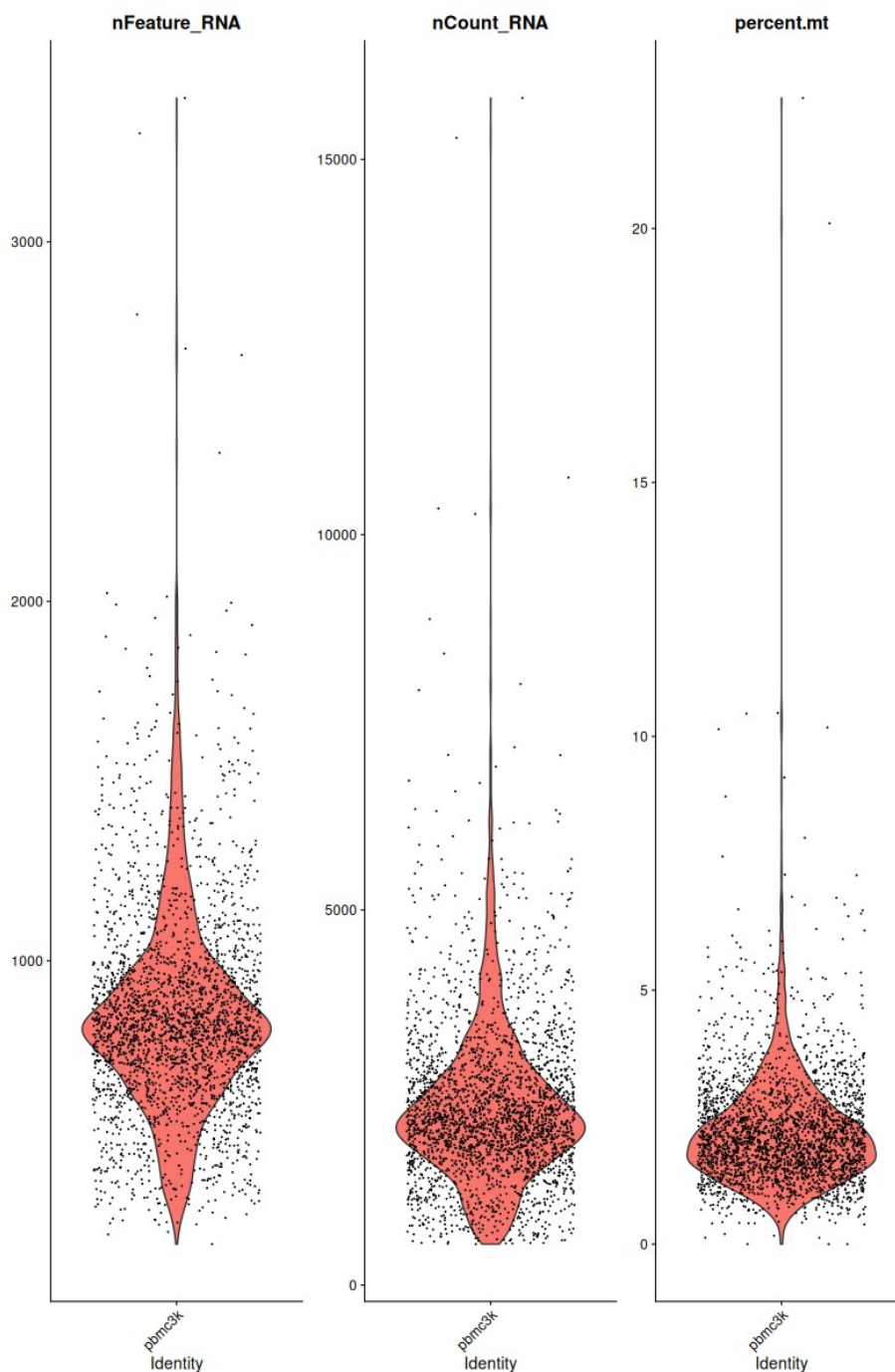


Figure 2: The scatter plots below illustrating the relationship between sequencing depth and mitochondrial percentage (left), and the correlation between sequencing depth and detected gene counts (right) used to assess cell quality prior to filtering.

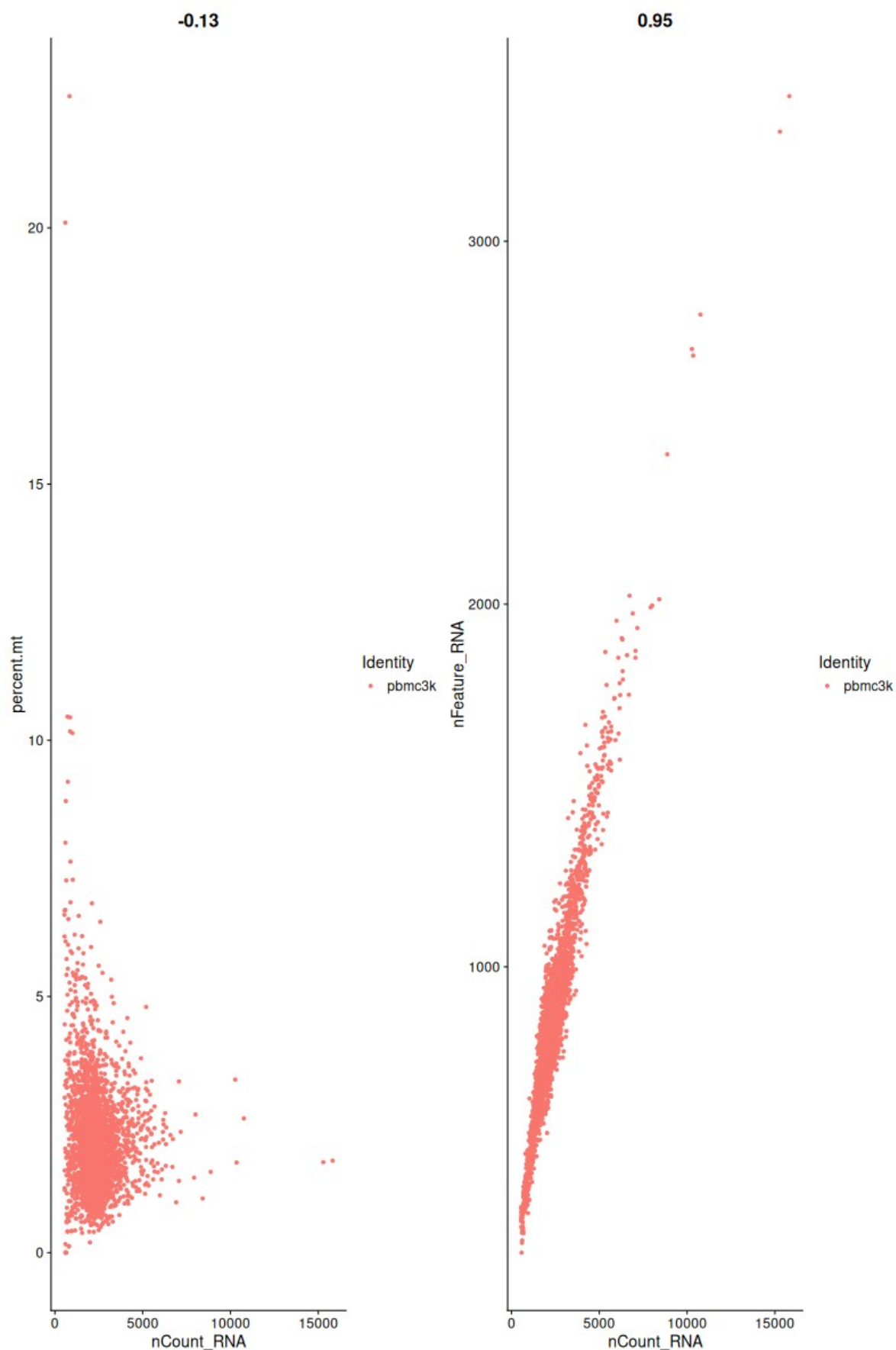


Figure 3: The scatter plots below displaying the relationship between average expression and standardized variance to identify highly variable features, with the top 10 most variable genes labeled in the right panel.

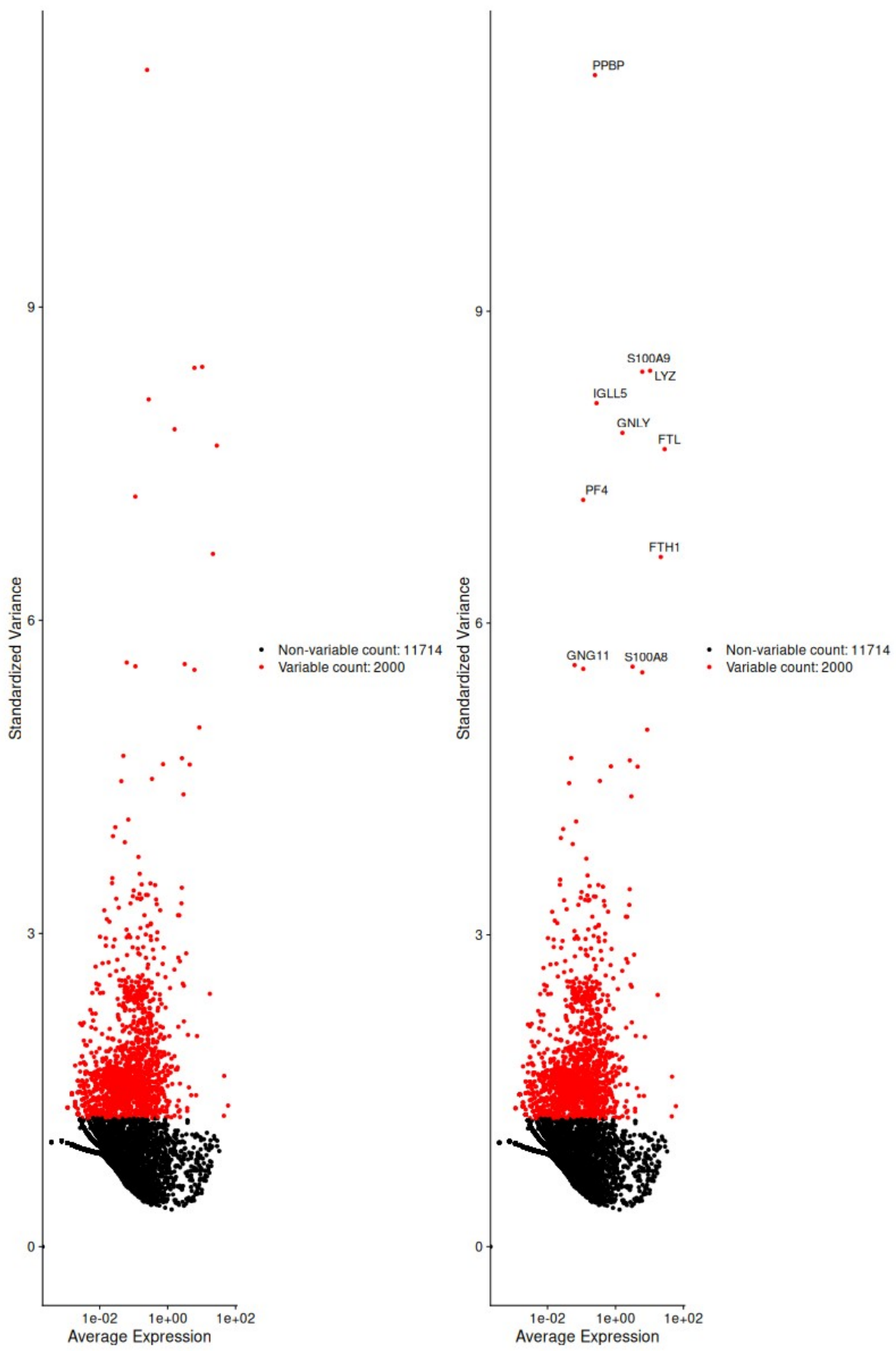


Figure 4: The plot displaying the top contributing genes for the first two principal components, ranked by their loading scores to illustrate the specific features responsible for the greatest variance in the dataset.

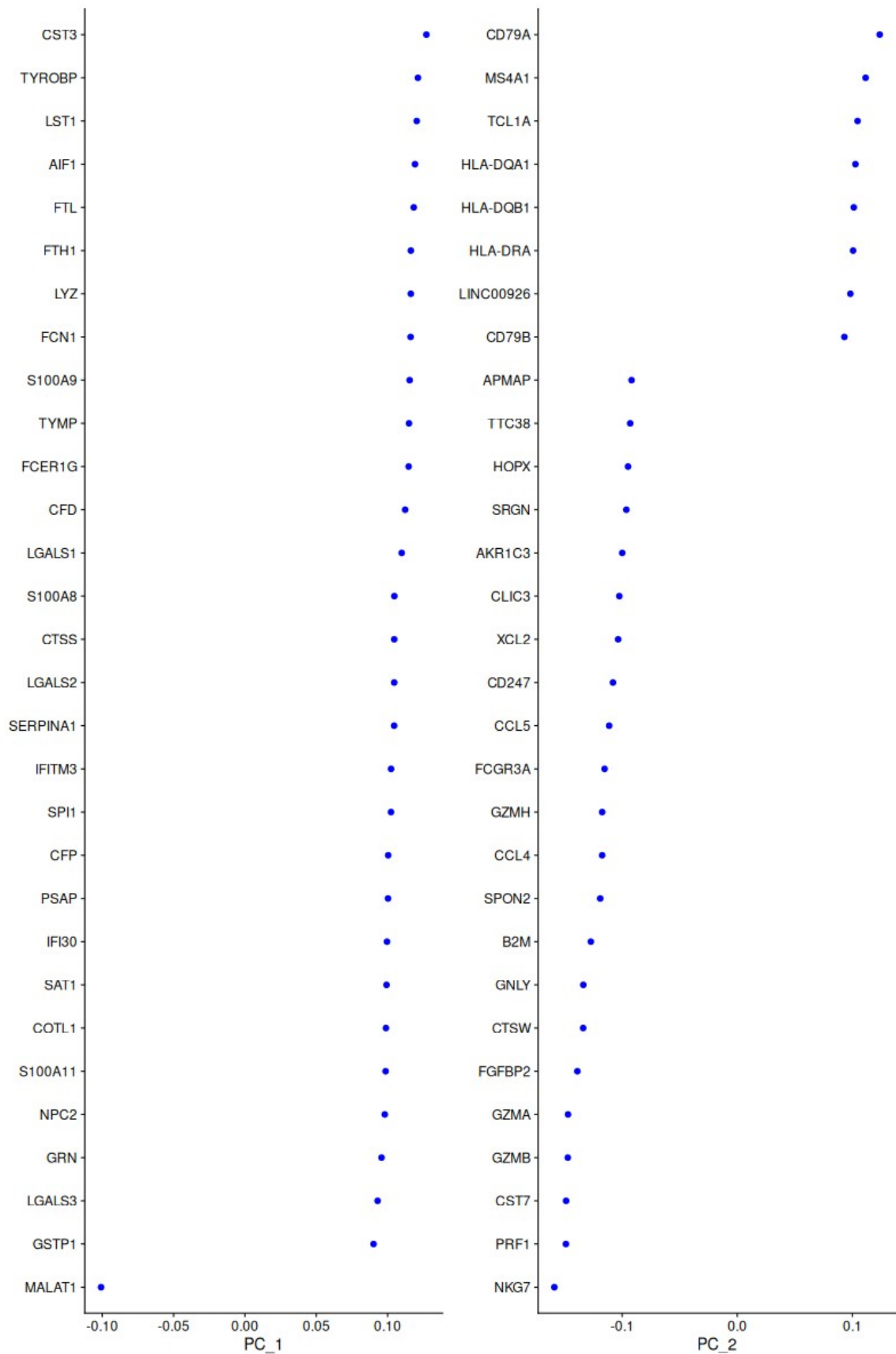


Figure 5: The plot visualizing the distribution of individual cells along the first two principal components (PC1 vs. PC2) of a Principal Component Analysis (PCA), to reveal global structure of our data distribution.

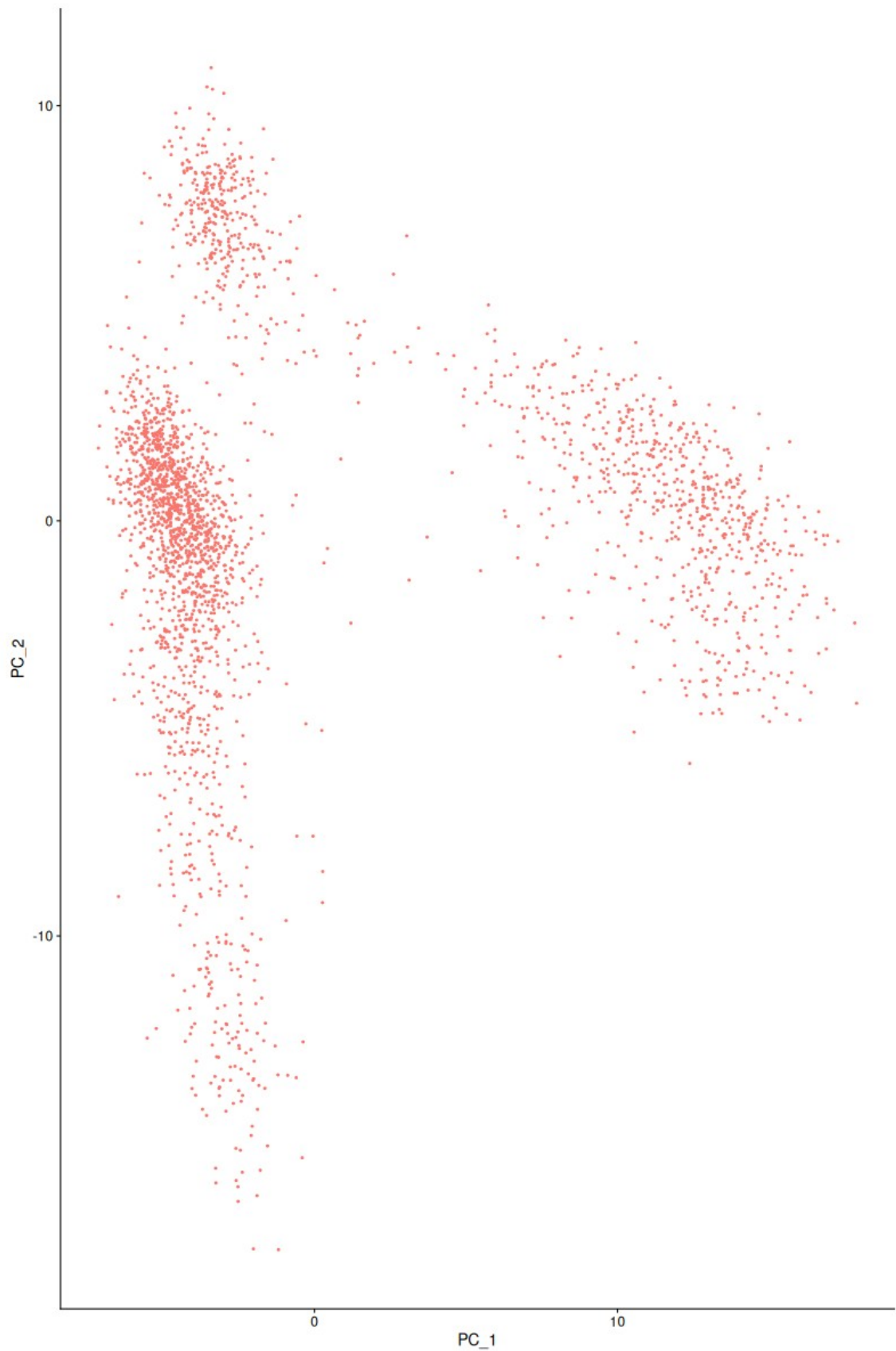


Figure 6: The following heatmap displaying the expression levels of the top contributing genes for the first principal component across 500 cells, highlighting the primary sources of heterogeneity in that dimension.

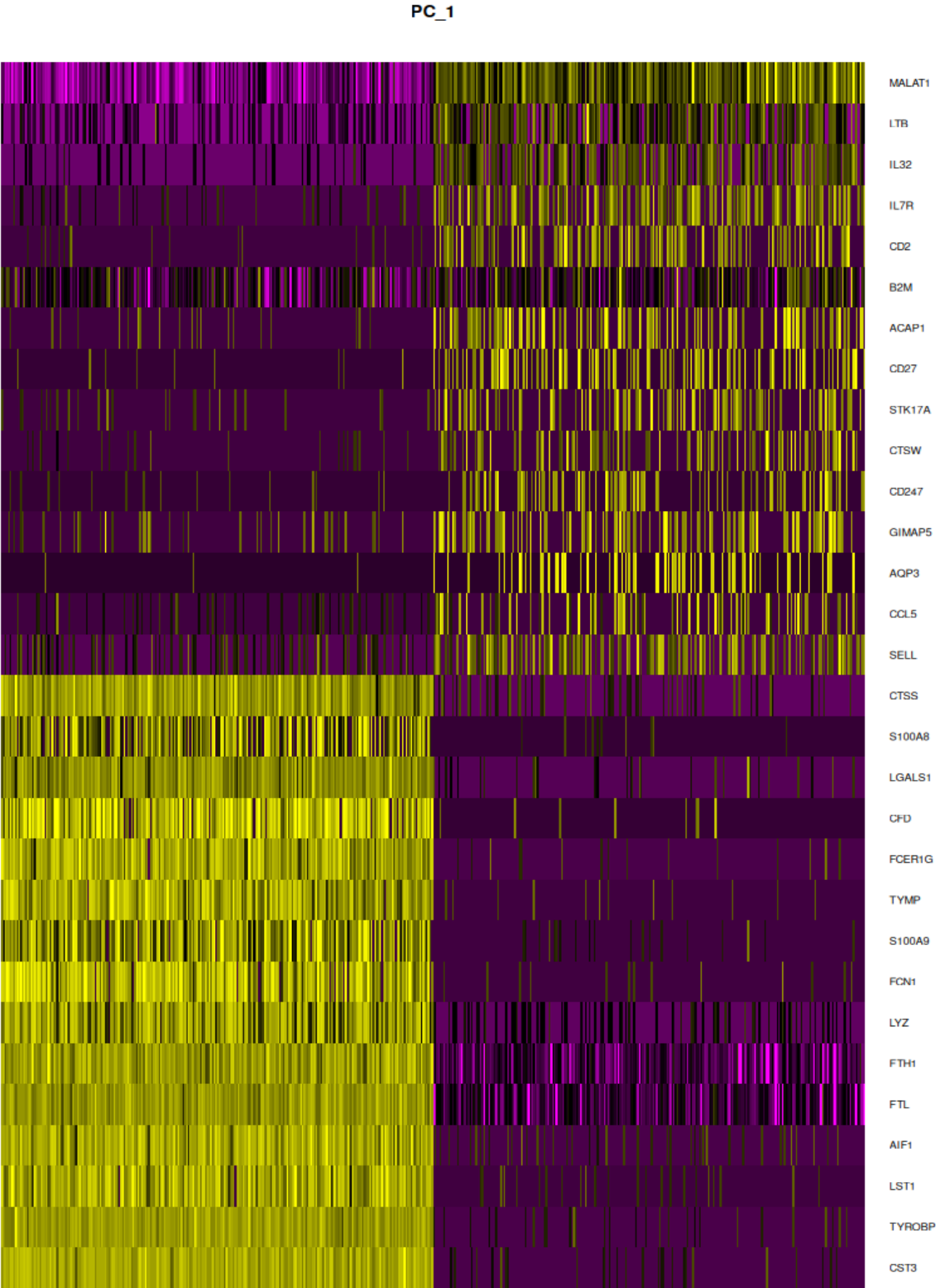


Figure 2 displays 15 heatmaps (PC_1 to PC_15) showing gene expression levels across 15 principal components. The genes are listed on the right of each heatmap, and the expression levels are color-coded: red for upregulated, green for downregulated, and black for non-significant.

PC_1: MAI AT1, LTB, IL1R1, CD4, ACAP1, CD37, STK17A, C15W, CD34, GIMAP5, AQP9, CD45, CD4, CD8, S100A8, GJA1, CD1, CCR10, TYMP, S100A9, CD4, LYZ, FTL, AIF1, CS1, TYROBP, CS12.

PC_2: NGG7, PRF1, CTS7, GZMB, GZMA, GCFP2, CTSW, GNL, CD4, CD34, GZMH, GZMB, CD34, CD4, CD37, HLA-DQB1, CD37, HLA-DQA2, HLA-DPA1, HLA-DQA, CD34, HLA-DPB1, CD37, UNCOR005, CD34, HLA-DQB1, HLA-DQA1, CD34, MSAA1, CD34.

PC_3: PRP, PRF, SOR, SPARC, GNG11, NKXIN, GPR, RGS18, TUBB1, CU, HST1H2AC, ARH1H9, TIGL, CD3, CD34, HLA-DQB1, UNCOR005, CD34, HLA-DQA2, HLA-DPA1, HLA-DQA, HLA-DPB1, HLA-DQB1, HLA-DQA1, CD34, CD37, HLA-DQA1, HLA-DQA1.

PC_4: VIM, IL1R1, S100A8, IL32, S100A9, S100A4, GIMAP7, S100A10, S100A9, MAI, AQP9, CD37, CD4, FVB, GJA1, S100A8, PRP, HLA-DQA2, HLA-DPA1, HLA-DQA1, CD34, SOR, PRF, HLA-DPB1, HST1H2AC, CD34, HLA-DQB1, MSAA1, CD34, CD37, HLA-DQA1.

PC_5: LTR, IL1R1, CD3, VIM, MSAA1, AQP9, CTSW, S100A8, GZMA, GCFP2, CTSW, GNL, CD4, CD37, HLA-DQB1, CD37, HLA-DQA2, HLA-DPA1, HLA-DQA, CD34, HLA-DPB1, CD37, UNCOR005, CD34, HLA-DQB1, HLA-DQA1, CD34, MSAA1, CD34.

PC_6: PRP, PRF, SOR, SPARC, GNG11, NKXIN, GPR, RGS18, TUBB1, CU, HST1H2AC, ARH1H9, TIGL, CD3, CD34, HLA-DQB1, UNCOR005, CD34, HLA-DQA2, HLA-DPA1, HLA-DQA, HLA-DPB1, HLA-DQB1, HLA-DQA1, CD34, CD37, HLA-DQA1, HLA-DQA1.

PC_7: GZMK, CD4, GZMA, LTR, S100A8, S100A4, GZMA, CD37, STK17A, C15W, CD34, GIMAP5, AQP9, CD45, CD4, CD8, S100A8, GJA1, CD1, CCR10, TYMP, S100A9, CD4, LYZ, FTL, AIF1, CS1, TYROBP, CS12.

PC_8: PRP, PRF, SOR, SPARC, GNG11, NKXIN, GPR, RGS18, TUBB1, CU, HST1H2AC, ARH1H9, TIGL, CD3, CD34, HLA-DQB1, UNCOR005, CD34, HLA-DQA2, HLA-DPA1, HLA-DQA, HLA-DPB1, HLA-DQB1, HLA-DQA1, CD34, CD37, HLA-DQA1, HLA-DQA1.

PC_9: PRP, PRF, SOR, SPARC, GNG11, NKXIN, GPR, RGS18, TUBB1, CU, HST1H2AC, ARH1H9, TIGL, CD3, CD34, HLA-DQB1, UNCOR005, CD34, HLA-DQA2, HLA-DPA1, HLA-DQA, HLA-DPB1, HLA-DQB1, HLA-DQA1, CD34, CD37, HLA-DQA1, HLA-DQA1.

PC_10: PRP, PRF, SOR, SPARC, GNG11, NKXIN, GPR, RGS18, TUBB1, CU, HST1H2AC, ARH1H9, TIGL, CD3, CD34, HLA-DQB1, UNCOR005, CD34, HLA-DQA2, HLA-DPA1, HLA-DQA, HLA-DPB1, HLA-DQB1, HLA-DQA1, CD34, CD37, HLA-DQA1, HLA-DQA1.

PC_11: PRP, PRF, SOR, SPARC, GNG11, NKXIN, GPR, RGS18, TUBB1, CU, HST1H2AC, ARH1H9, TIGL, CD3, CD34, HLA-DQB1, UNCOR005, CD34, HLA-DQA2, HLA-DPA1, HLA-DQA, HLA-DPB1, HLA-DQB1, HLA-DQA1, CD34, CD37, HLA-DQA1, HLA-DQA1.

PC_12: PRP, PRF, SOR, SPARC, GNG11, NKXIN, GPR, RGS18, TUBB1, CU, HST1H2AC, ARH1H9, TIGL, CD3, CD34, HLA-DQB1, UNCOR005, CD34, HLA-DQA2, HLA-DPA1, HLA-DQA, HLA-DPB1, HLA-DQB1, HLA-DQA1, CD34, CD37, HLA-DQA1, HLA-DQA1.

PC_13: PRP, PRF, SOR, SPARC, GNG11, NKXIN, GPR, RGS18, TUBB1, CU, HST1H2AC, ARH1H9, TIGL, CD3, CD34, HLA-DQB1, UNCOR005, CD34, HLA-DQA2, HLA-DPA1, HLA-DQA, HLA-DPB1, HLA-DQB1, HLA-DQA1, CD34, CD37, HLA-DQA1, HLA-DQA1.

PC_14: PRP, PRF, SOR, SPARC, GNG11, NKXIN, GPR, RGS18, TUBB1, CU, HST1H2AC, ARH1H9, TIGL, CD3, CD34, HLA-DQB1, UNCOR005, CD34, HLA-DQA2, HLA-DPA1, HLA-DQA, HLA-DPB1, HLA-DQB1, HLA-DQA1, CD34, CD37, HLA-DQA1, HLA-DQA1.

PC_15: PRP, PRF, SOR, SPARC, GNG11, NKXIN, GPR, RGS18, TUBB1, CU, HST1H2AC, ARH1H9, TIGL, CD3, CD34, HLA-DQB1, UNCOR005, CD34, HLA-DQA2, HLA-DPA1, HLA-DQA, HLA-DPB1, HLA-DQB1, HLA-DQA1, CD34, CD37, HLA-DQA1, HLA-DQA1.

Figure 8: This elbow plot displaying the standard deviation of each principal component, used to visually identify the inflection point where true biological signal diminishes into random noise.

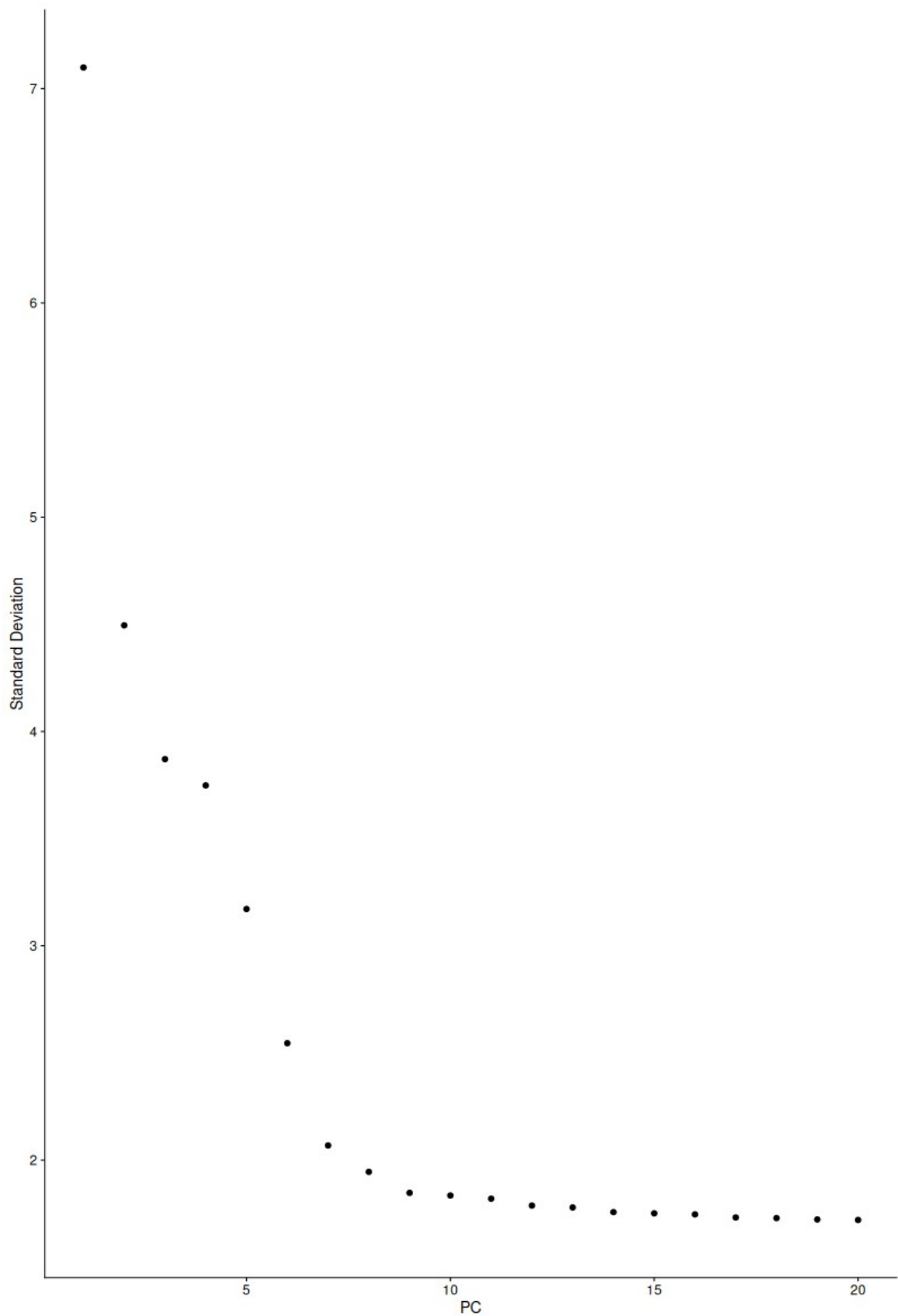


Figure 9: This shows a UMAP projection visualizing the single-cell clustering results, where each point represents an individual cell and colors indicate distinct clusters identified based on the first 10 principal components, with resolution 0.5.

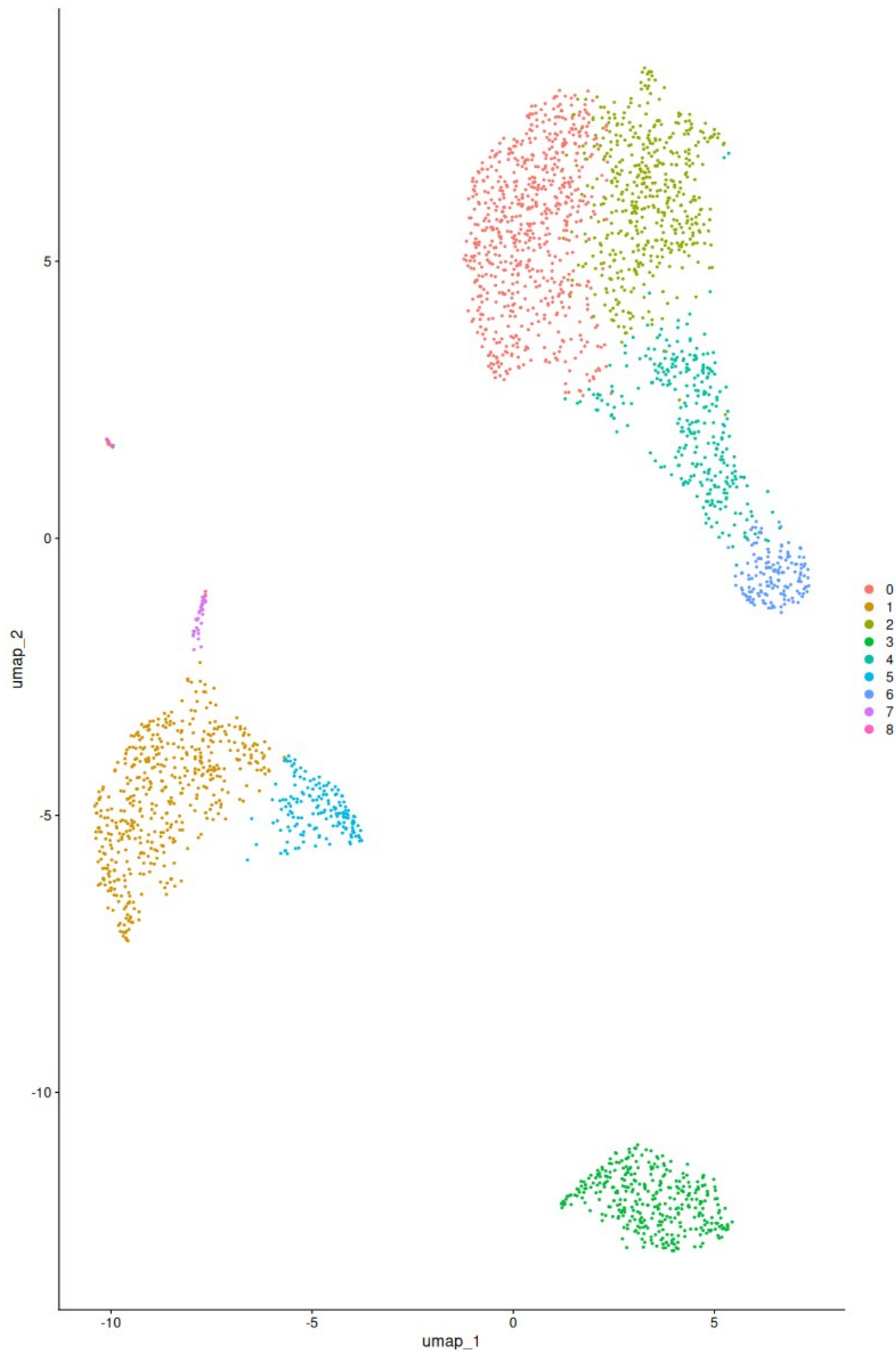


Figure 10: This shows a UMAP projection with the distinct clusters identified based on the first 5 principal components, with resolution 0.5. We observe that this affect our results and giving us a less clear clustering, loosing some clusters and the samples shows a broader distribution in the space compare with the figure 9.

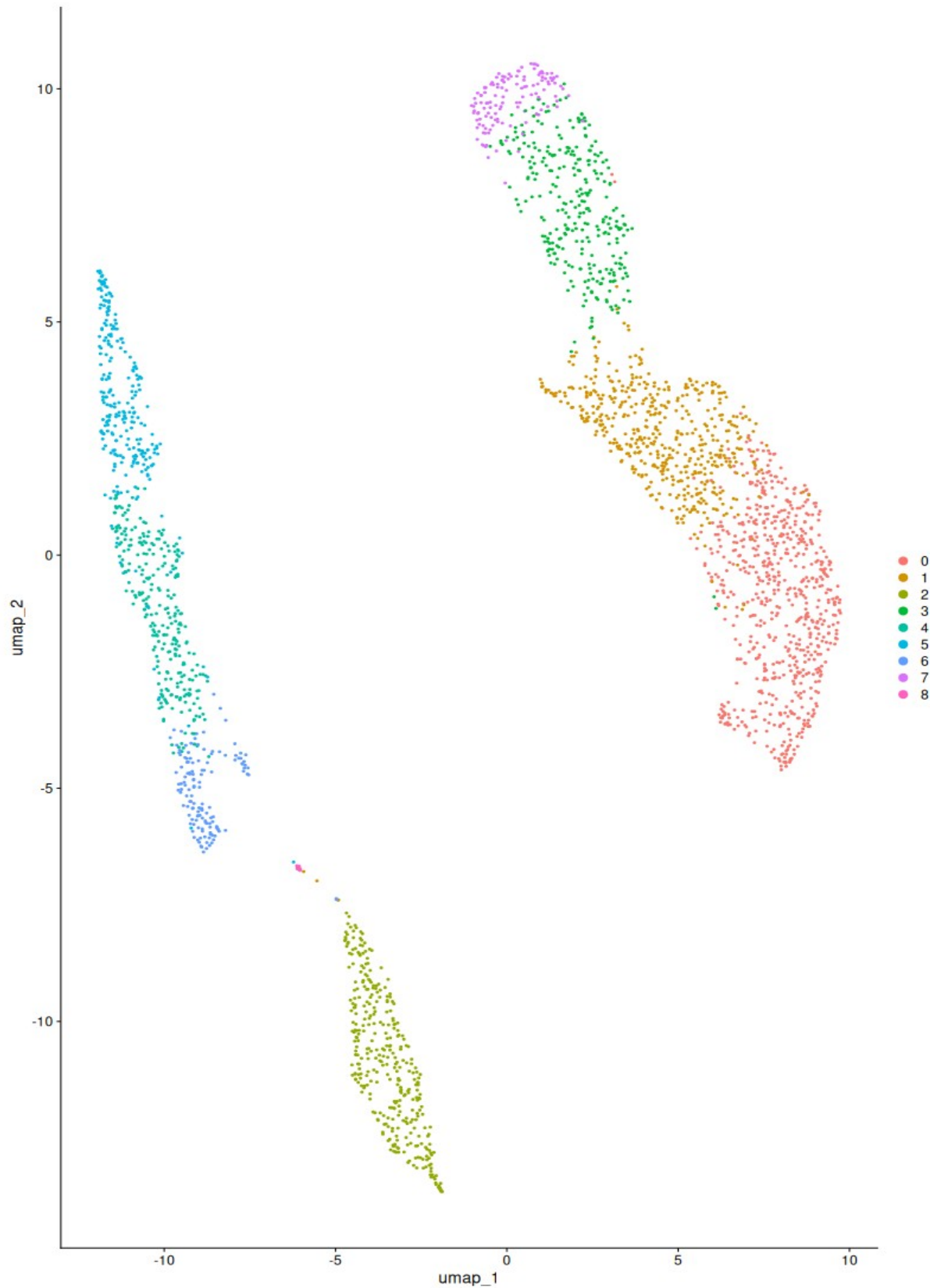


Figure 11: In this UMAP projection, we have used the same resolution (0.5), but we keep 15 principal components in the PCA. We observe that this affect the plot significant compare to figure 9 and 10 and show as a better separation among the clusters.

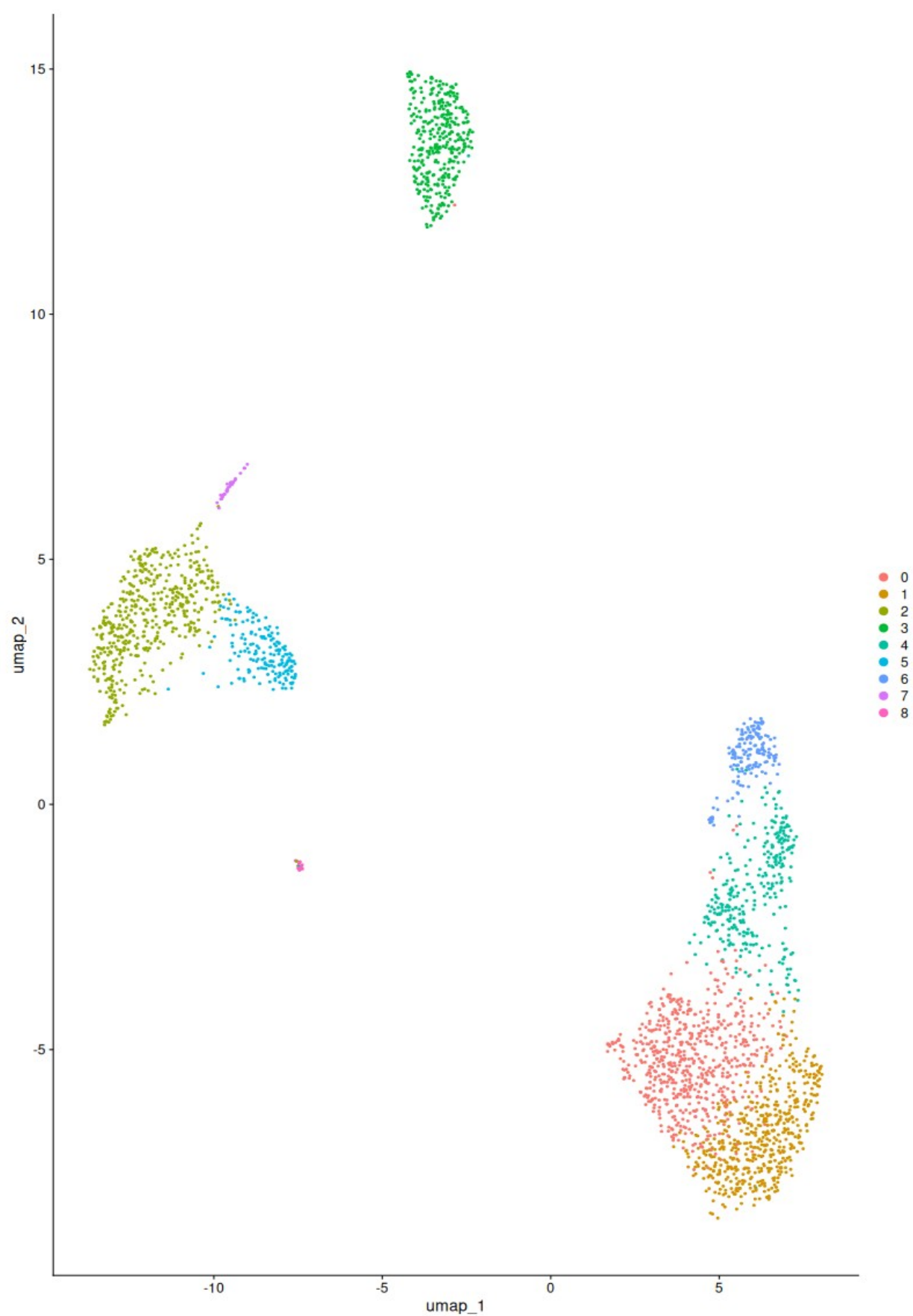


Figure 12: In this UMAP projection, we have maintain 10 principal components and change the resolution to 0.3. The position of the samples is the same as figure 9, but the number of the clusters looks decline like we zoom out and separate the samples to bigger subgroups, by merging the clusters that have a lot of similarities.

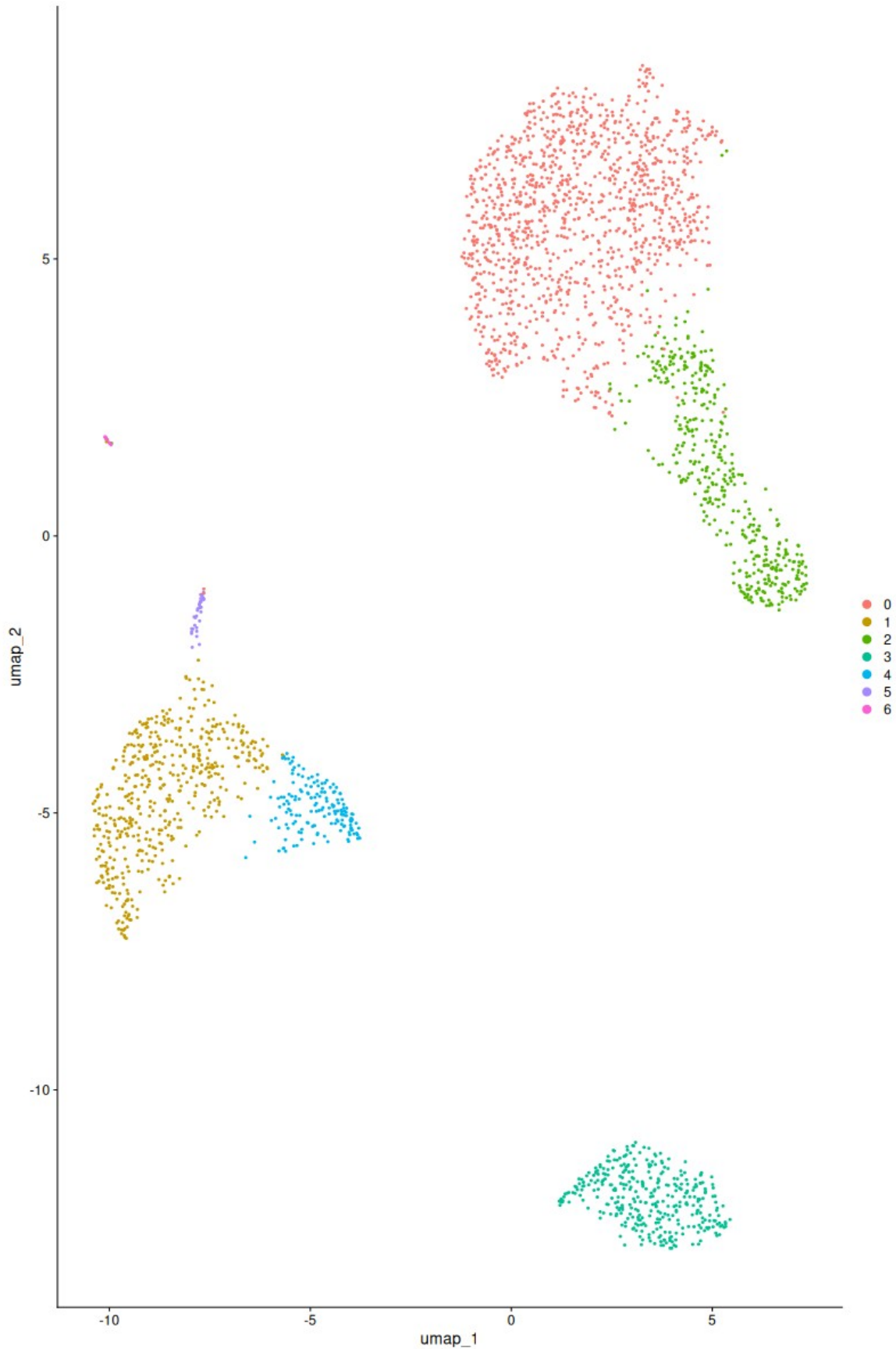


Figure 13: In this UMAP projection, we have maintain 10 principal components and change the resolution to 0.8. Here we observe that the position of the sample is maintain from figure 9, but contrary to figure 12, we are doing a zoom in and increasing the number of clusters, which can mean that we separate different subtypes of the same cell types.

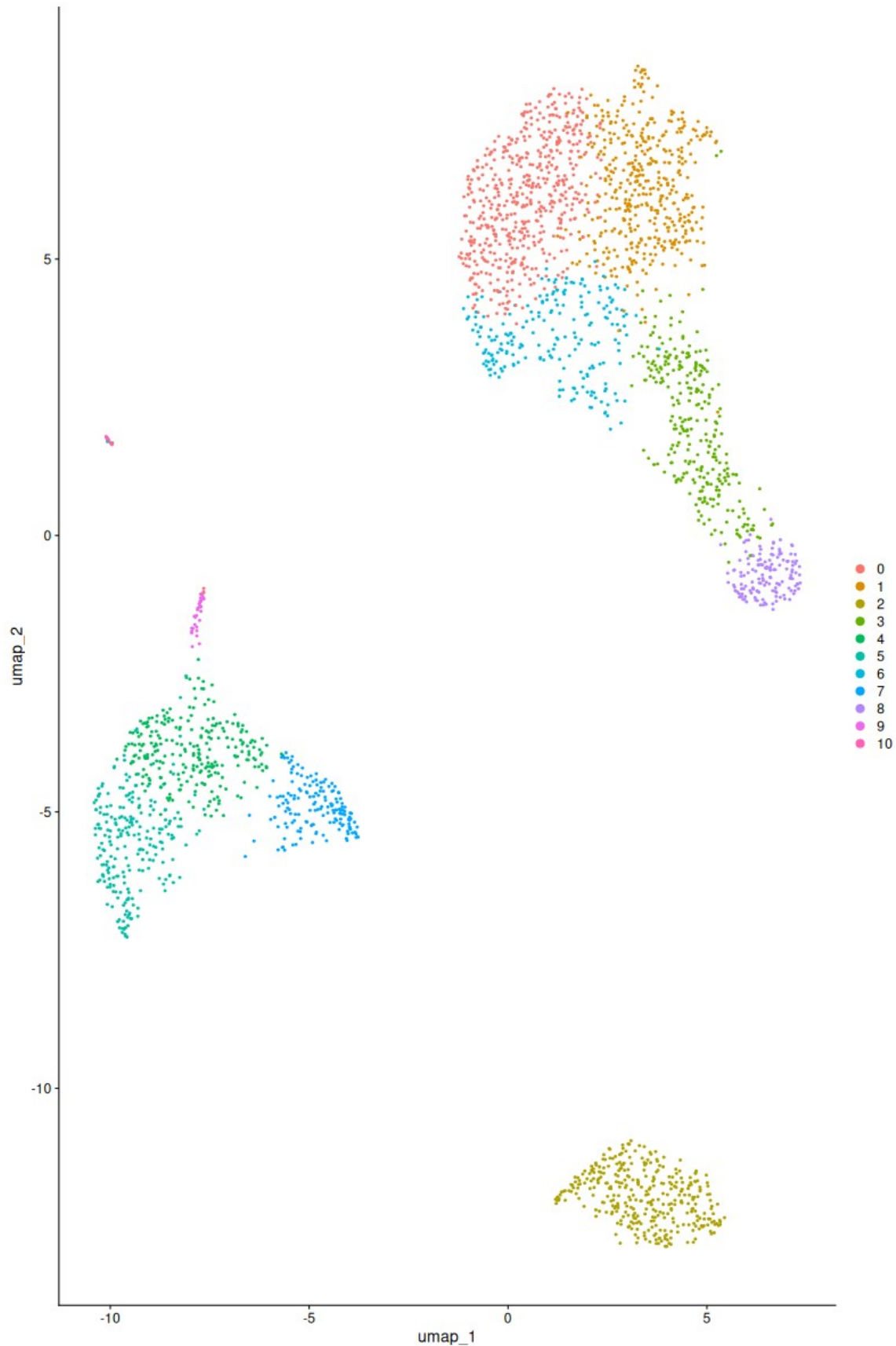


Figure 14: In this Violin plots is displayed the expression probability distributions of the B-cell lineage markers *MS4A1* and *CD79A* across all clusters, demonstrating their specific enrichment within a distinct cell population, which we expect to correspond to the B-cells population.

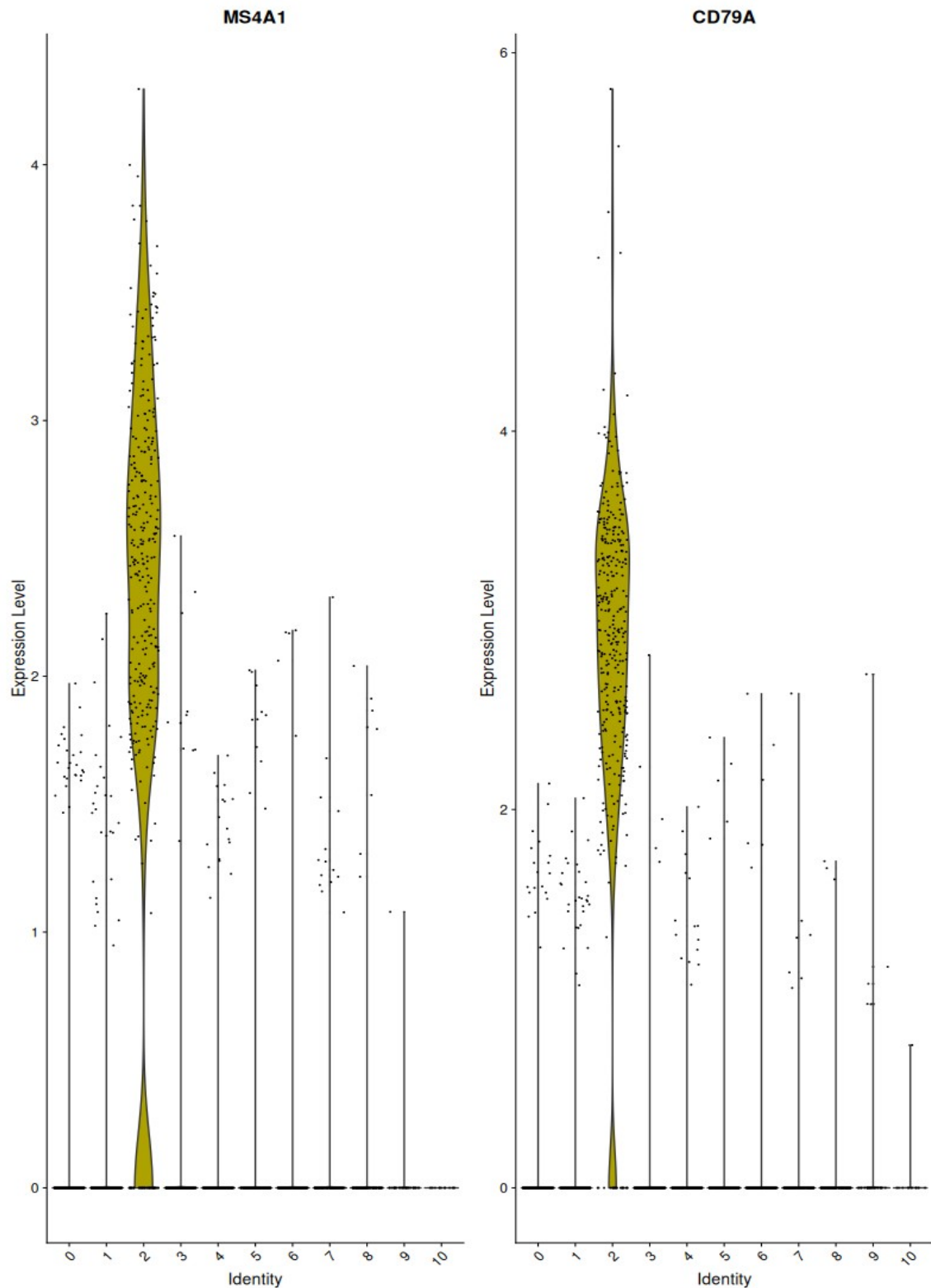


Figure 15: This violin plot displaying the distribution of log-transformed raw read counts for the markers *NKG7* and *PF4*, allowing for the assessment of absolute expression levels across clusters independent of normalization scaling.

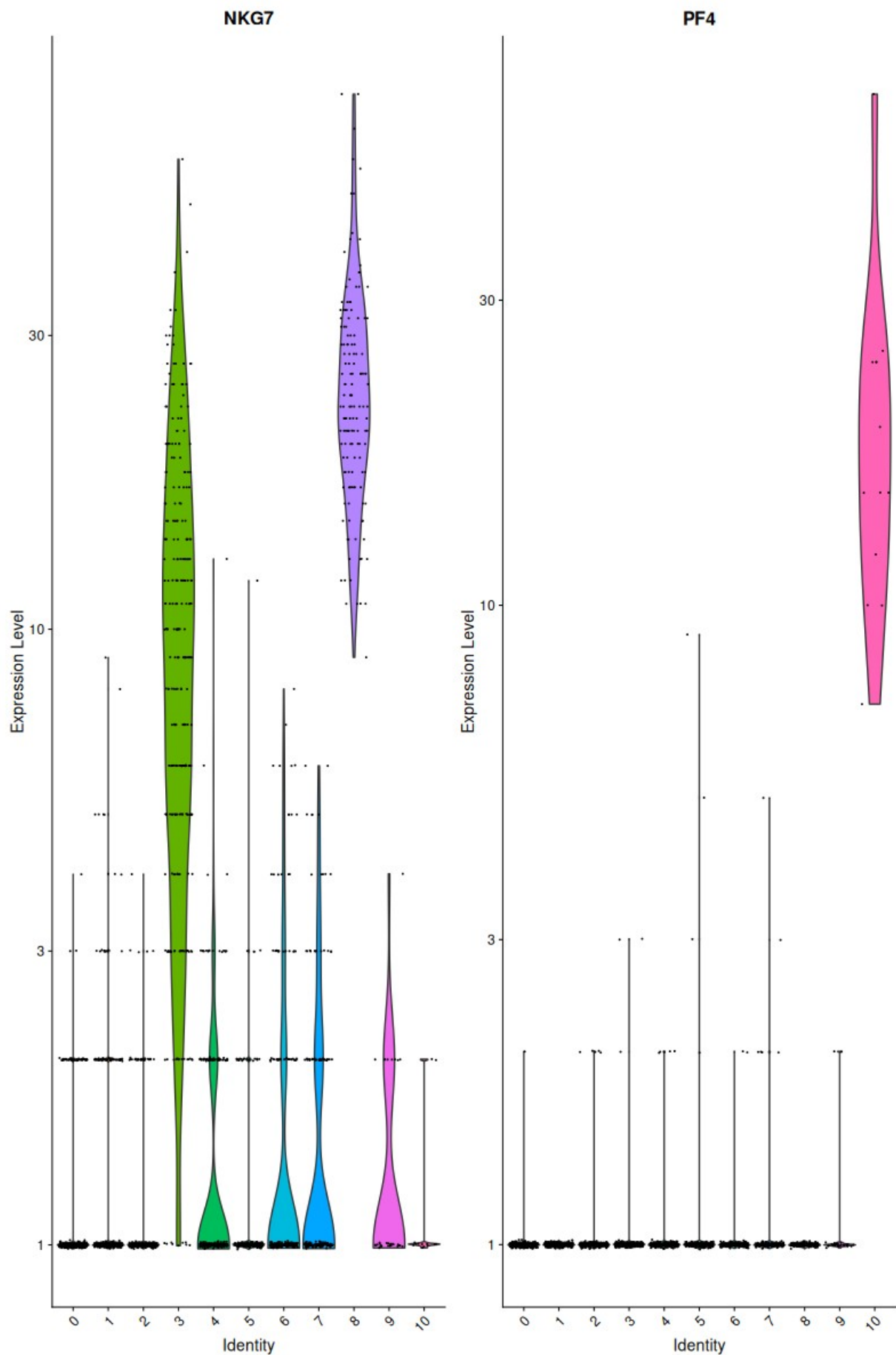


Figure 16: This figure plots visualizing the expression distribution of key canonical marker genes on the UMAP projection, used to map specific biological identities to the defined cell clusters.

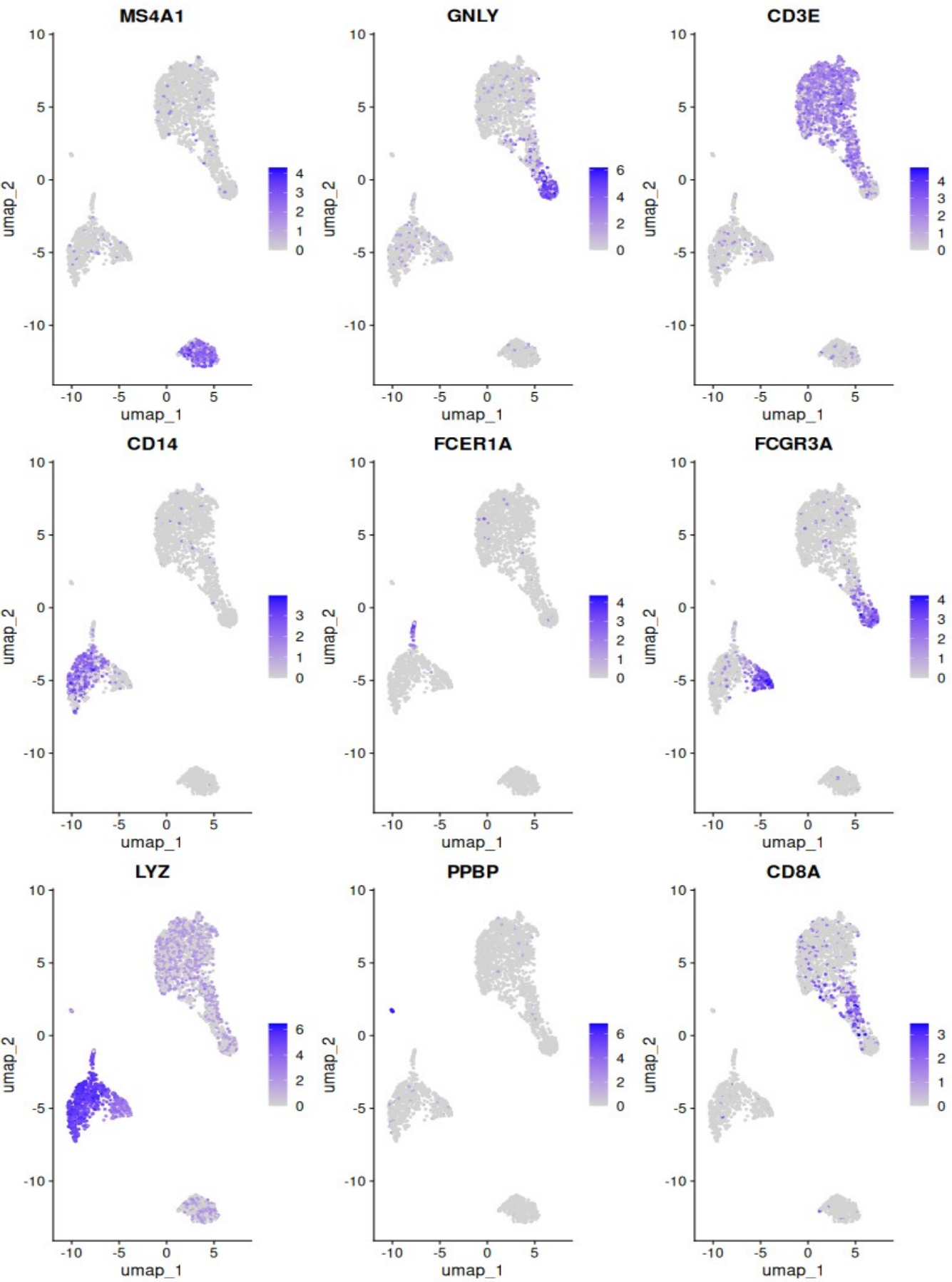


Figure 17: This figure plots illustrating the differential expression of *CCL5*, *MS4A1*, and *SEPT5* across the UMAP projection. We can observe that there is a distinct localization of Cytotoxic T/NK cells (*SEPT5*) and B cells (*MS4A1*), respectively, while Platelets (*SEPT5*) are not appear at all, which is consistent with our analysis, given the fact that we select only nuclear cell, during this experiment.

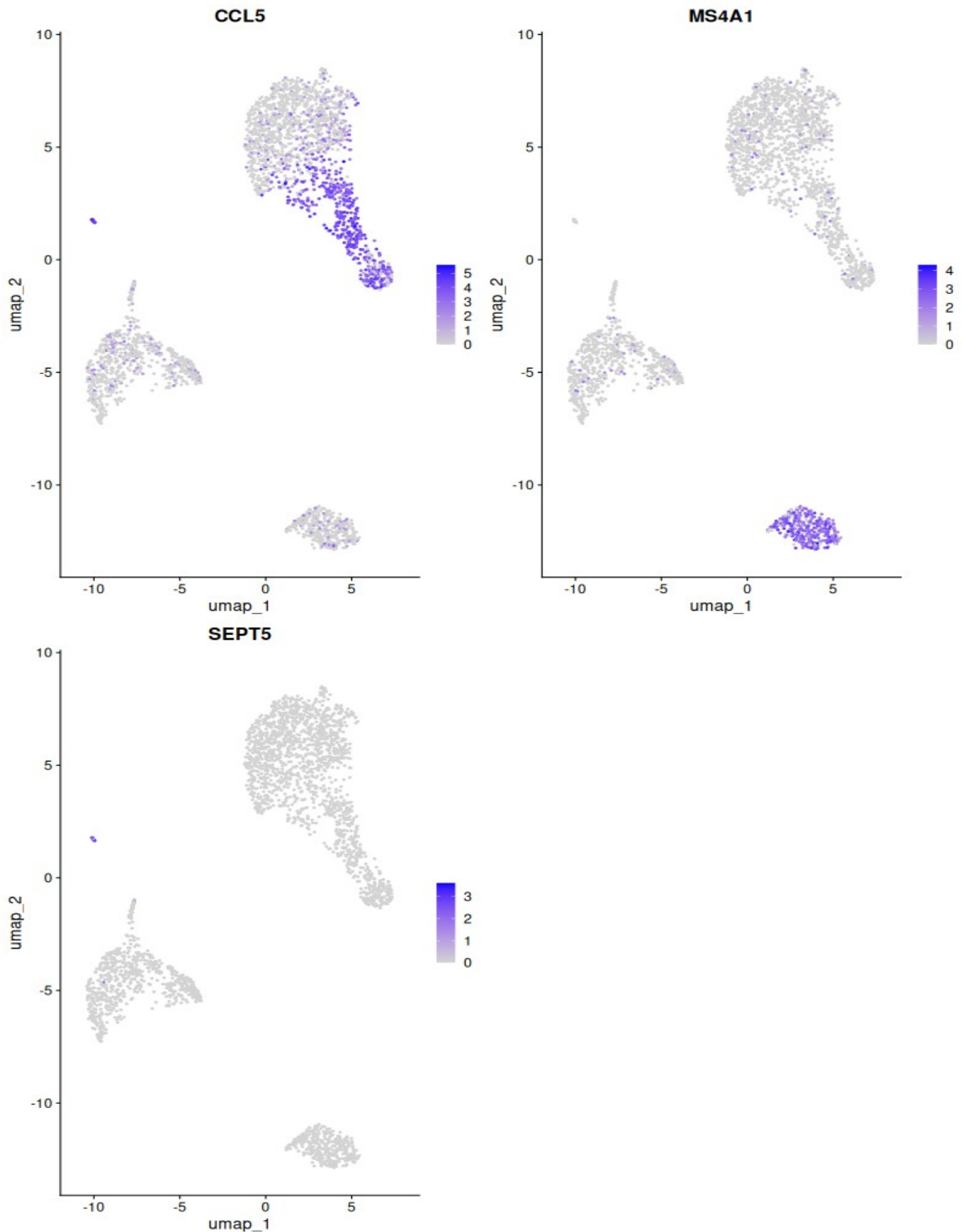


Figure 18: This heatmap visualizing the scaled expression of the top 10 upregulated marker genes for each cluster, displaying the distinct transcriptional signatures that define each cell population.

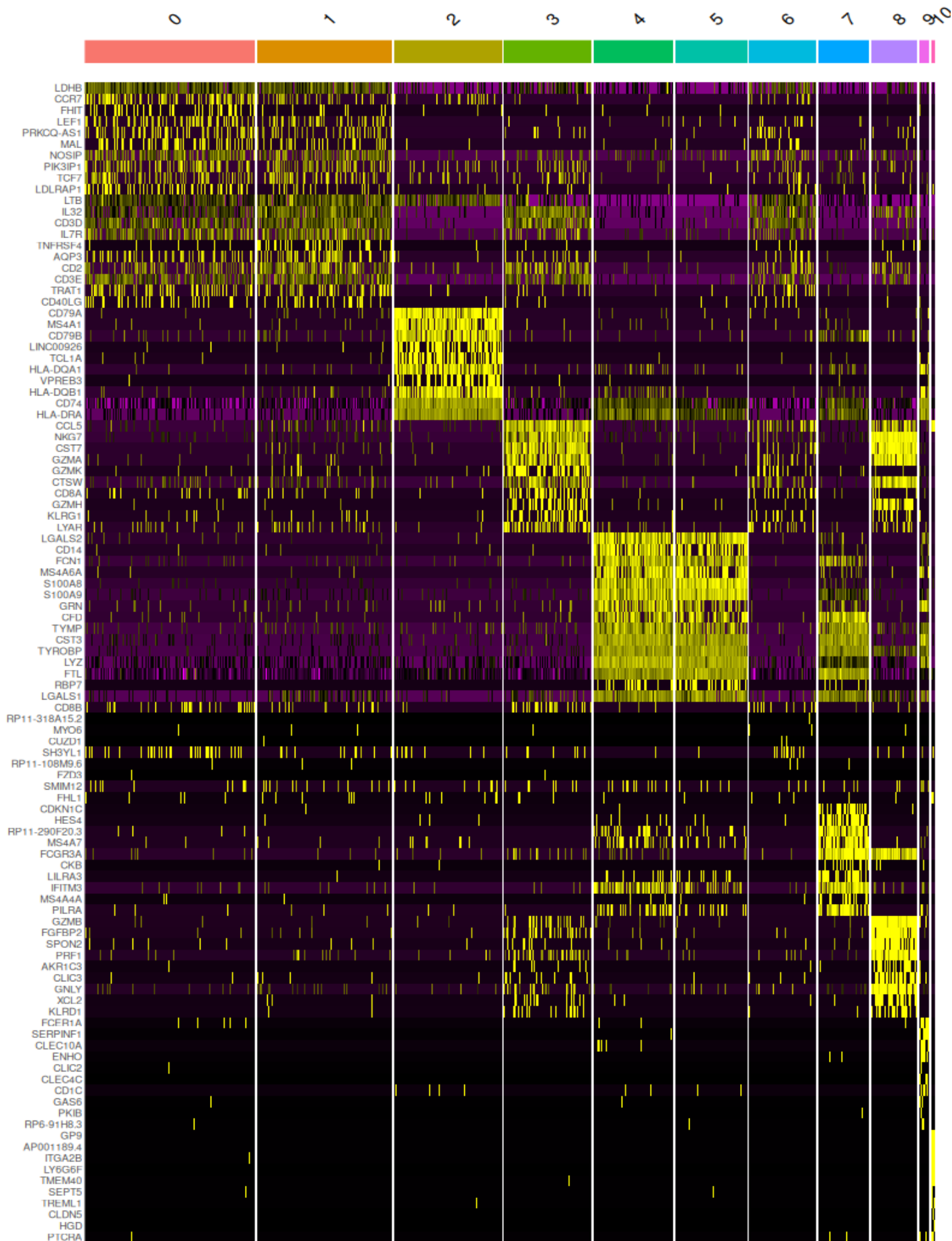


Figure 19: This UMAP projection showing the final cell type annotations, where the numerical cluster IDs have been replaced with biological labels derived from the canonical marker gene analysis.

