

Identifying related gene sets in colorectal cancer

Single Cell Sequencing Analytics Summer School 2022

Workflow pipeline

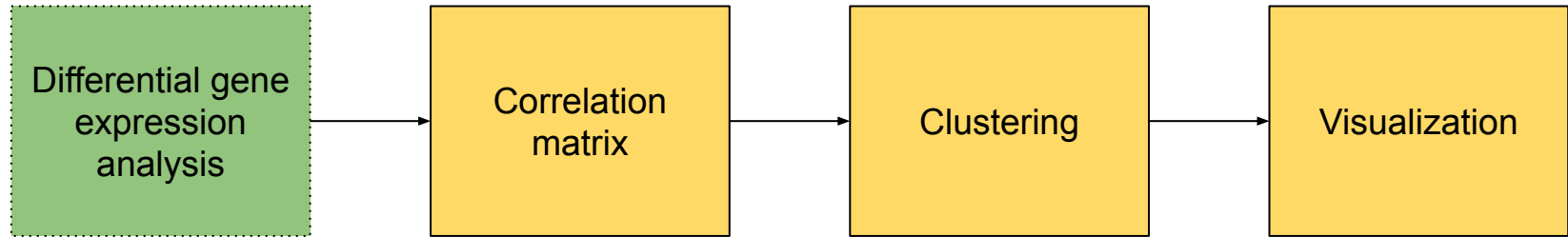
Colorectal cancer data set: cells from patients with and without tumor.



Question: Do the correlation patterns of genes in certain cell subtypes change in healthy vs. cancer tissue?

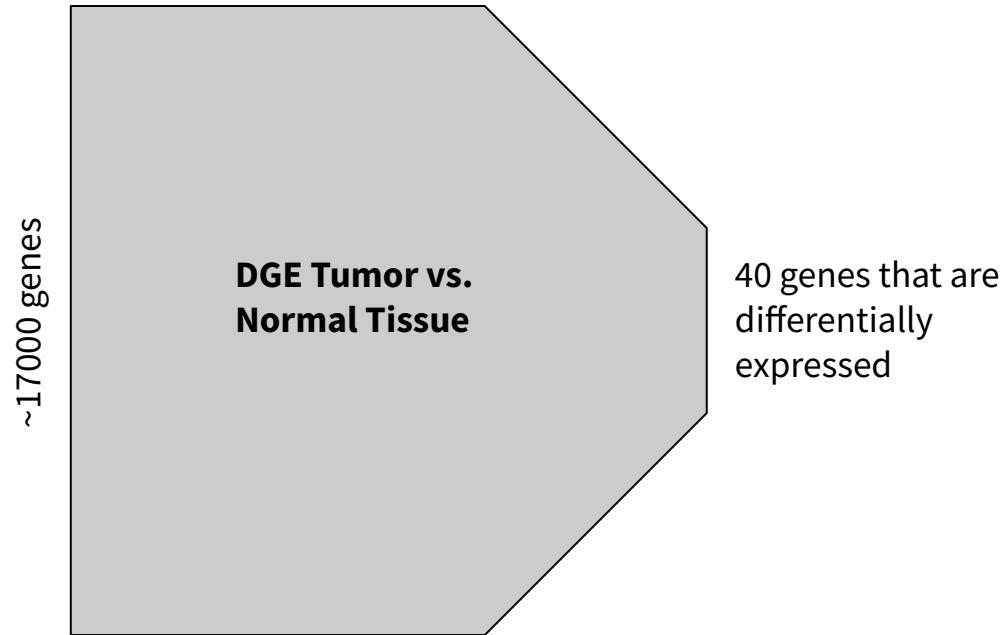
Workflow pipeline

Colorectal cancer data set: cells from patients with and without tumor.

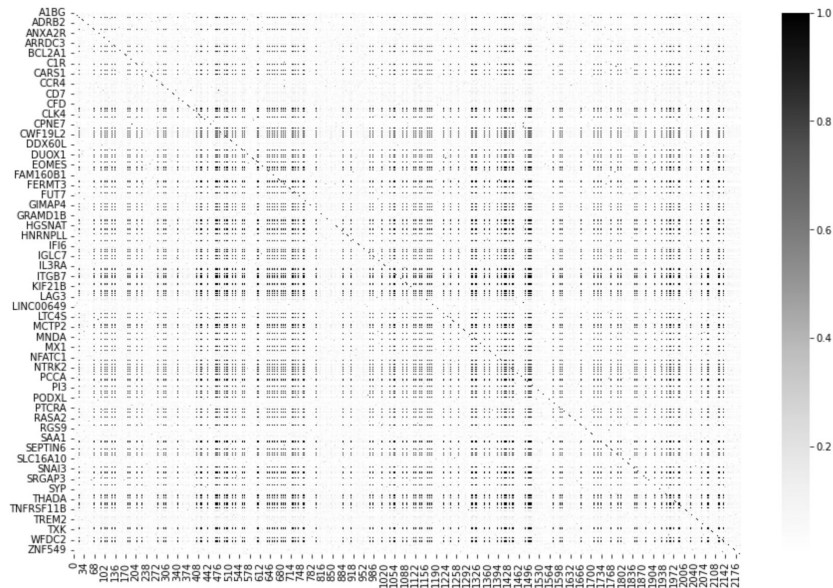


Question: Do the correlation patterns of genes in certain cell subtypes change in healthy vs. cancer tissue?

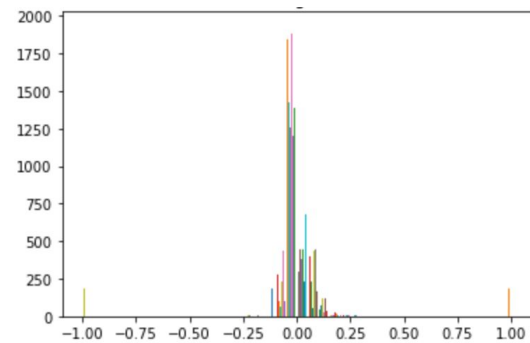
Differential gene expression analysis



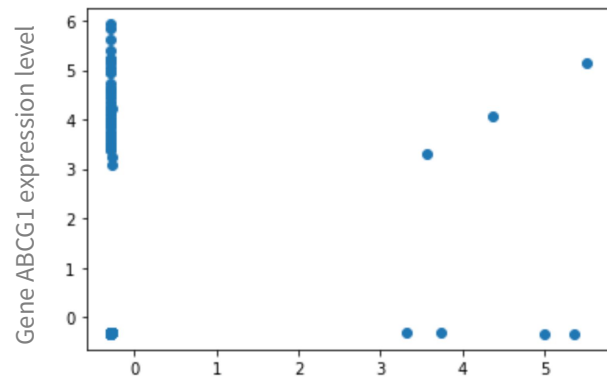
Pearson correlation



Heatmap of Pearson correlation matrix

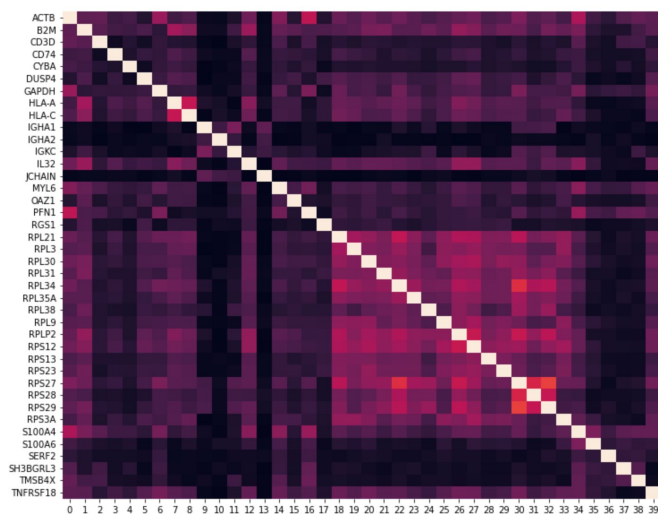


Histogram of the correlations



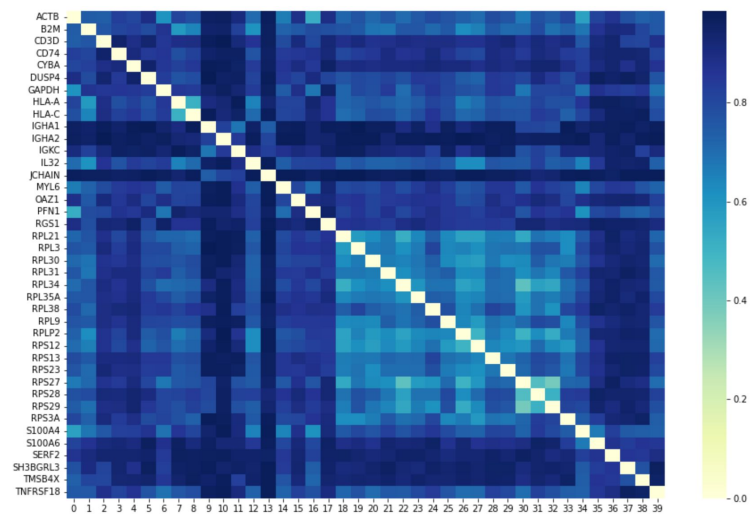
Gene ABCA1 expression level

Distance correlation (Székely et al., 2007)



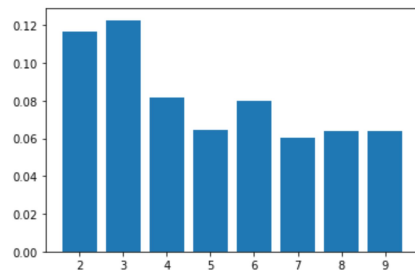
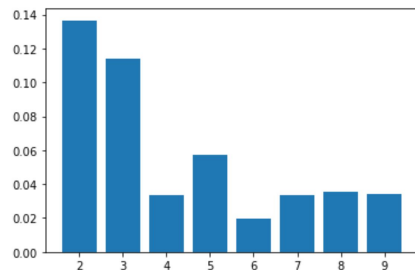
DC matrix

$$d(x,y) = 1 - |\rho(x,y)|$$

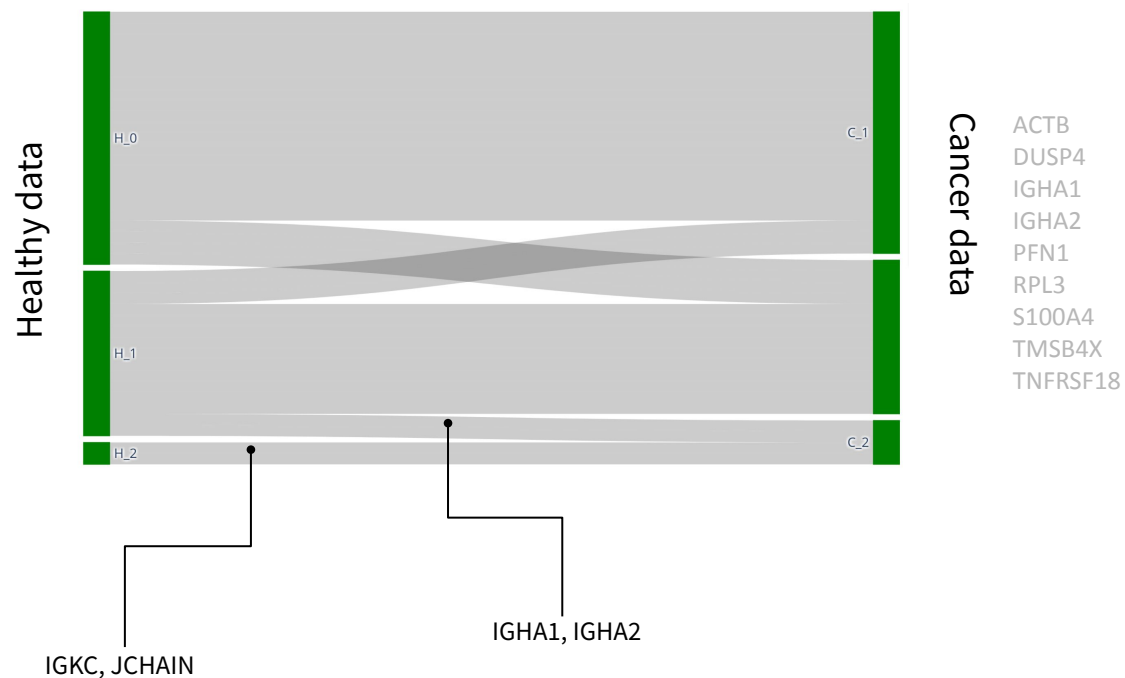


Distance matrix

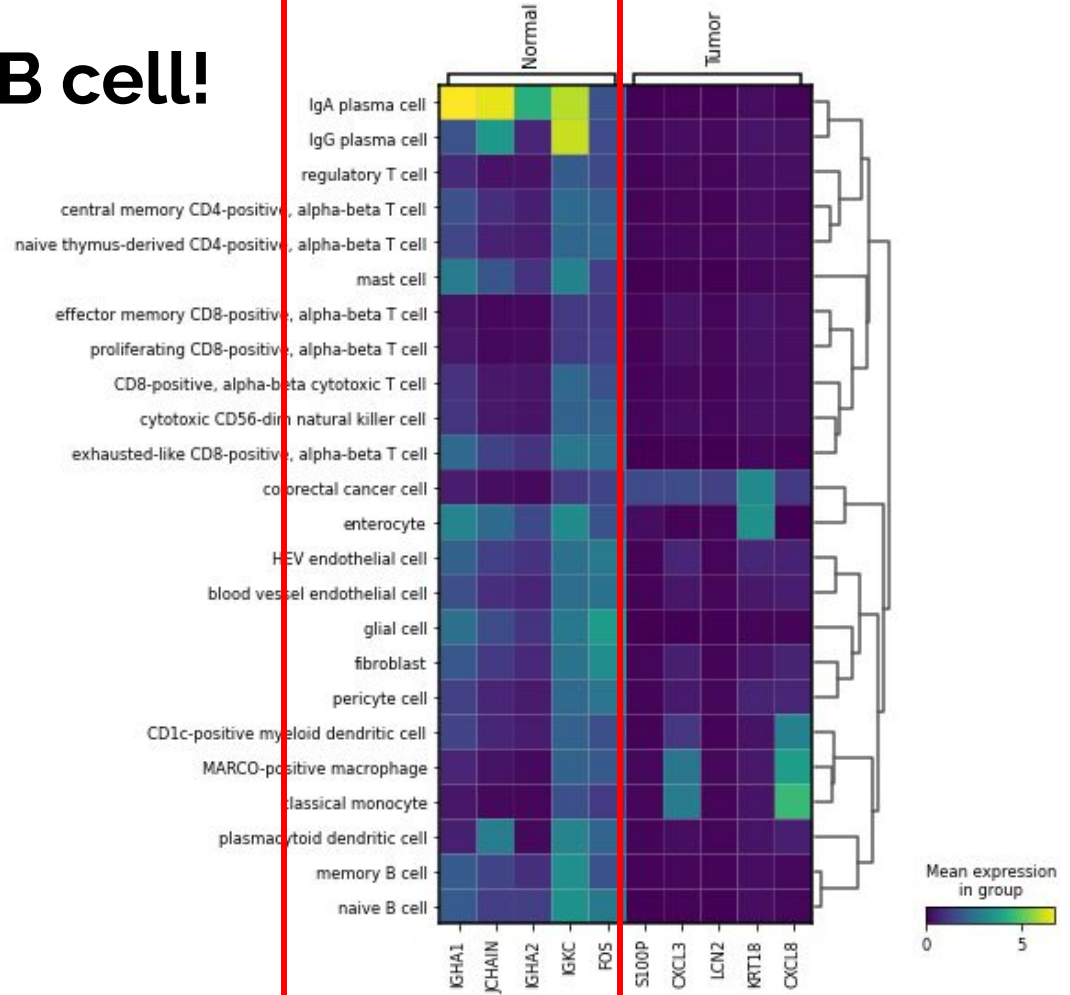
Comparing Clusters



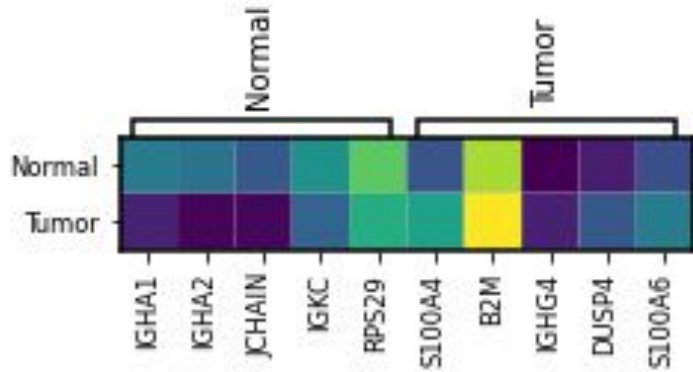
Average silhouette scores for healthy (top) and cancer (bottom) data



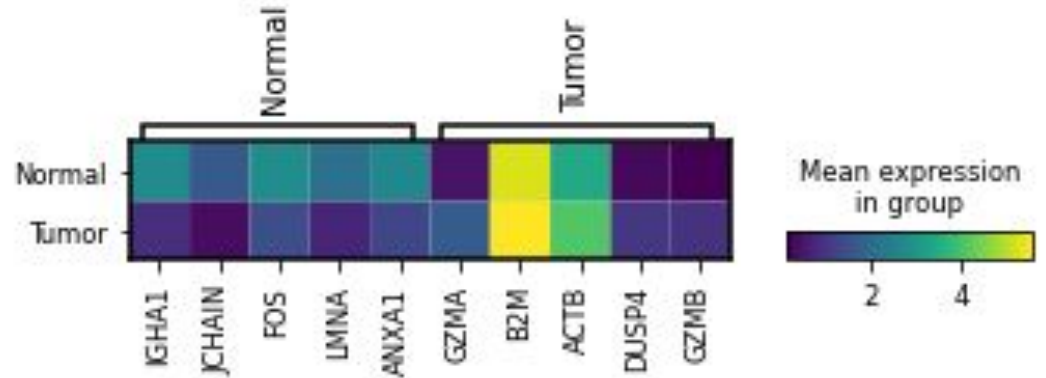
DGE - Everything is a B cell!



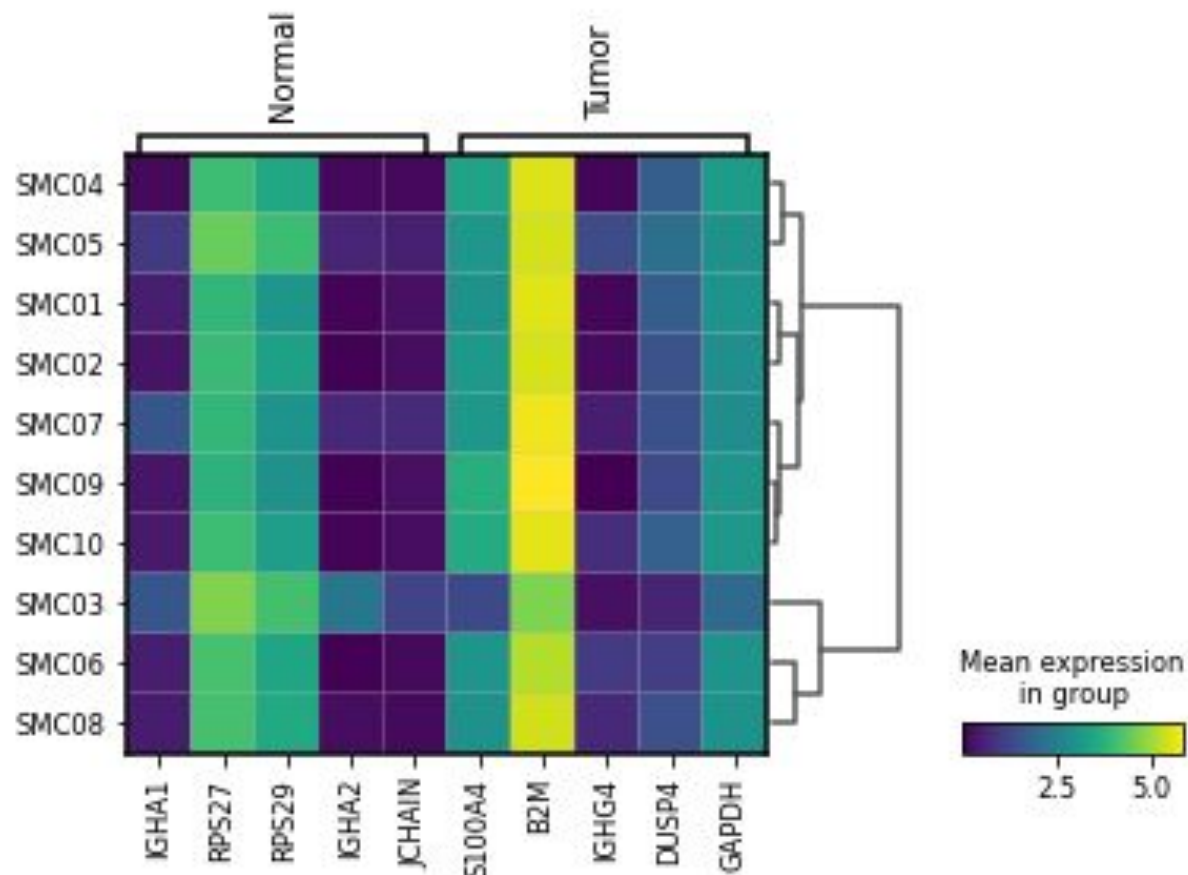
T-Reg DGE



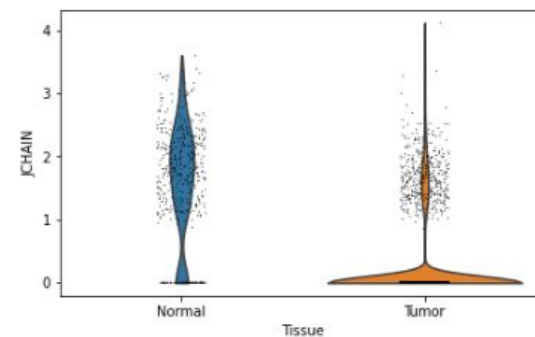
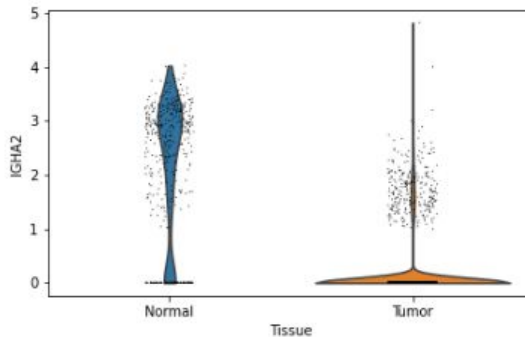
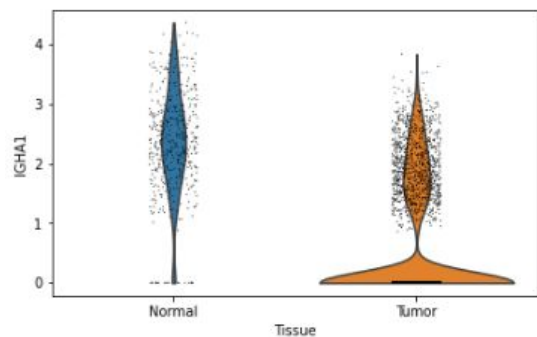
CD4 αβ+ DGE



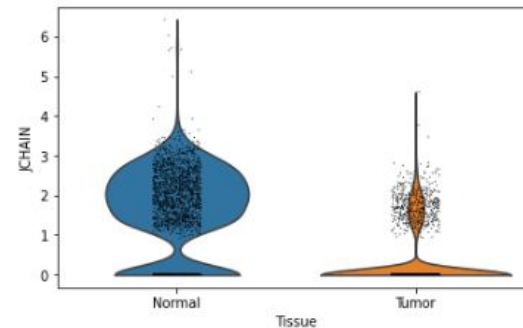
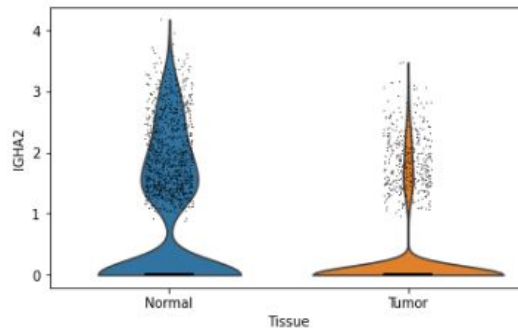
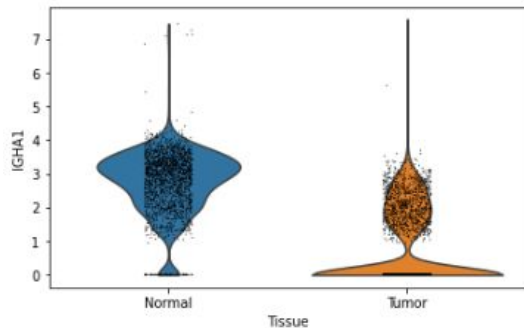
T-Reg DGE Per Patient



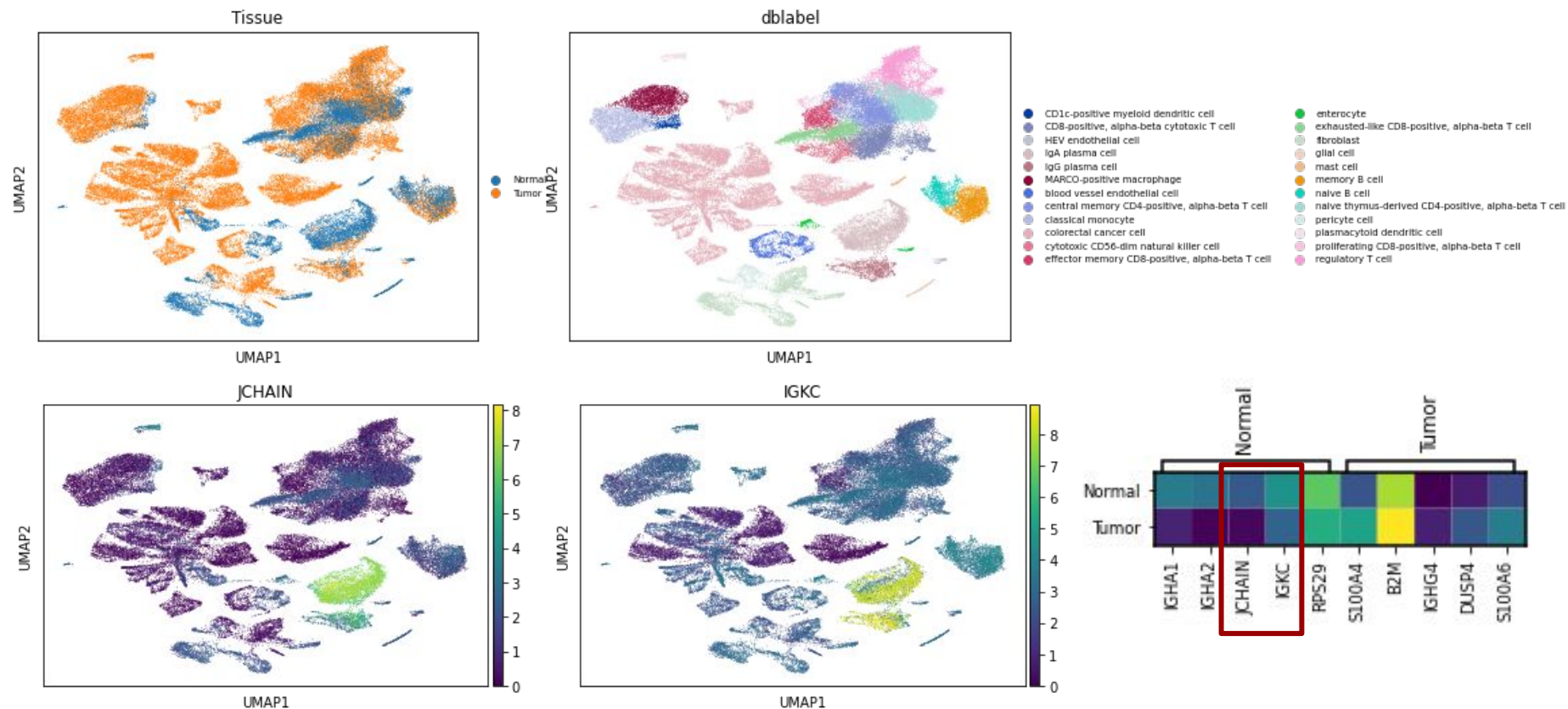
B cell gene expression among T-Regs



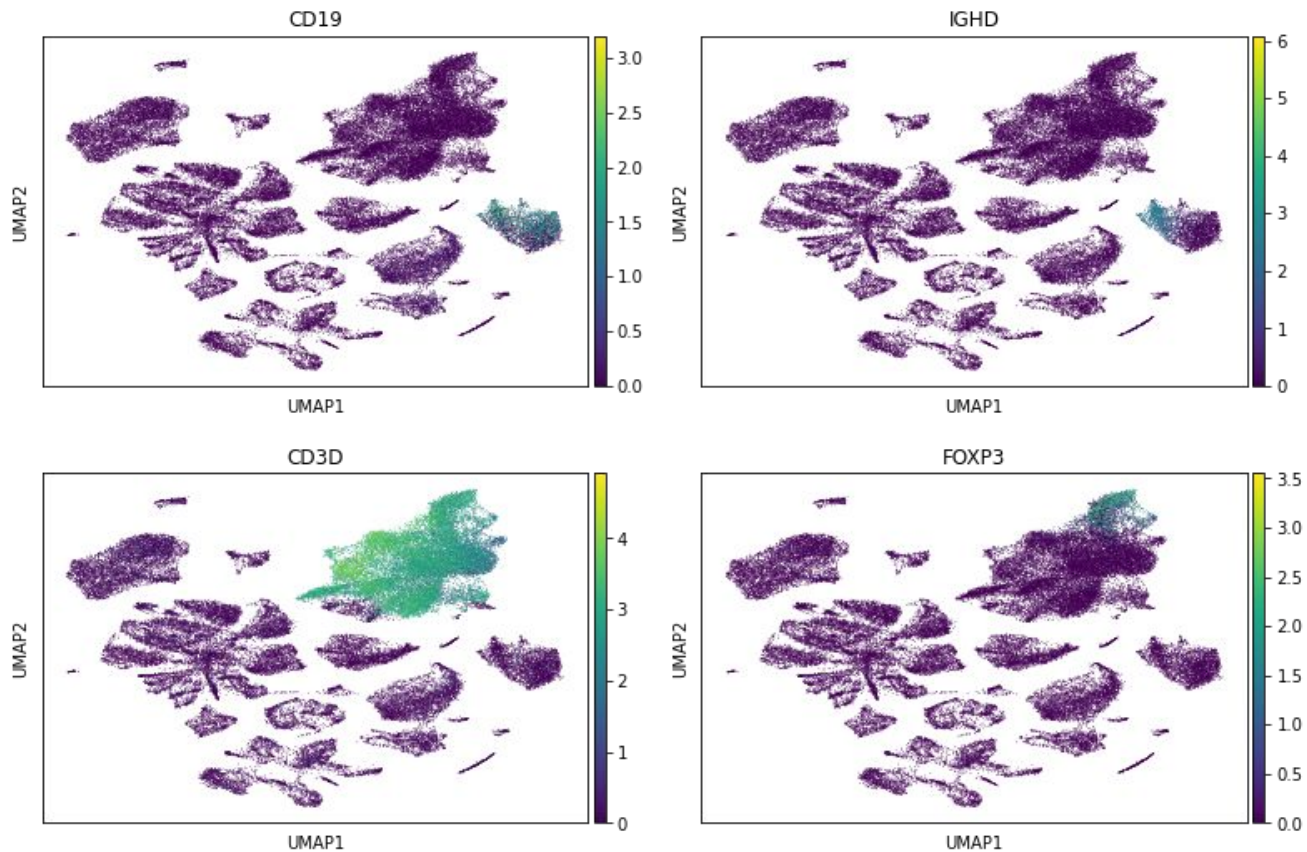
CD4+ $\alpha\beta$



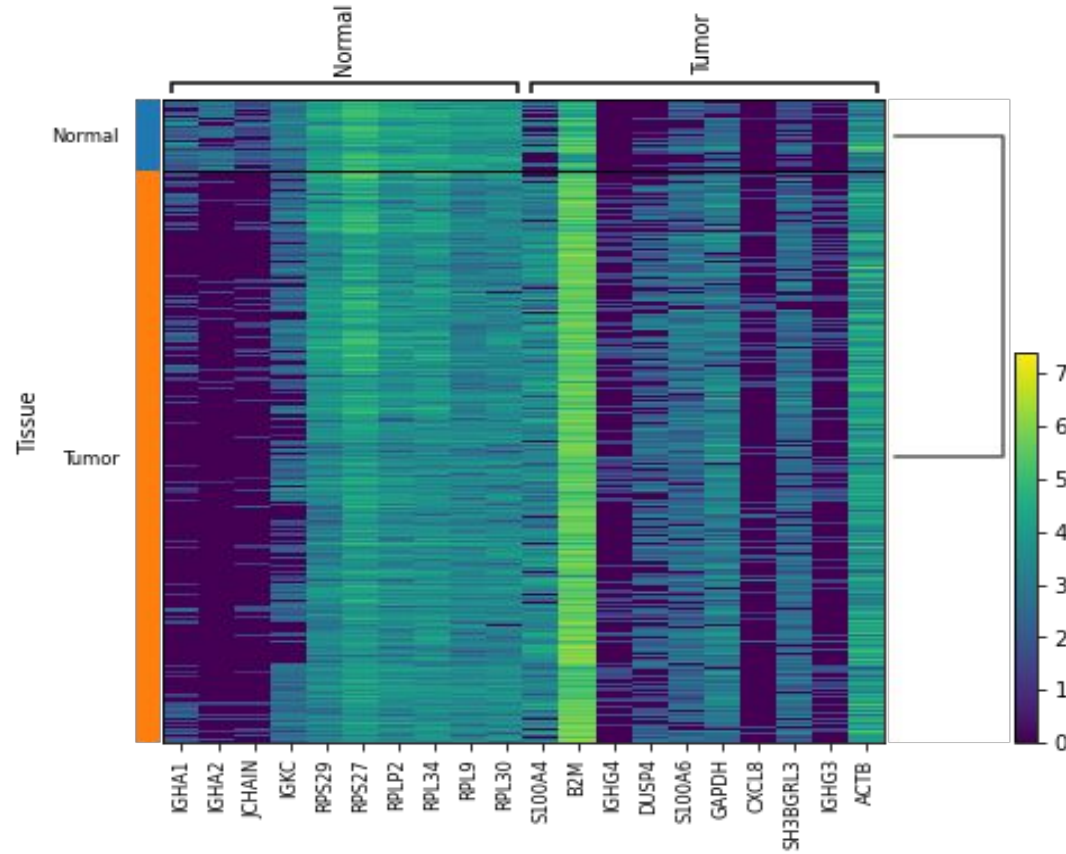
B cell gene expression



Not data is erroneous

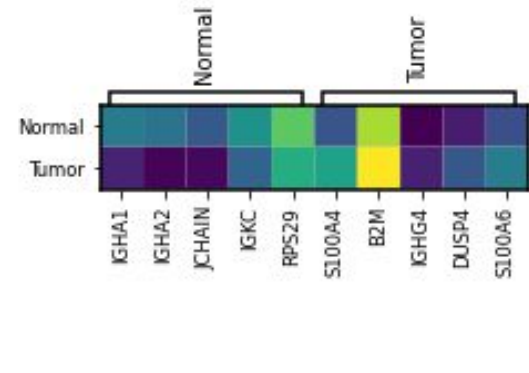


Unbalanced datasets



Number of Tregs in Tissue

Normal	Tumor
463	3716



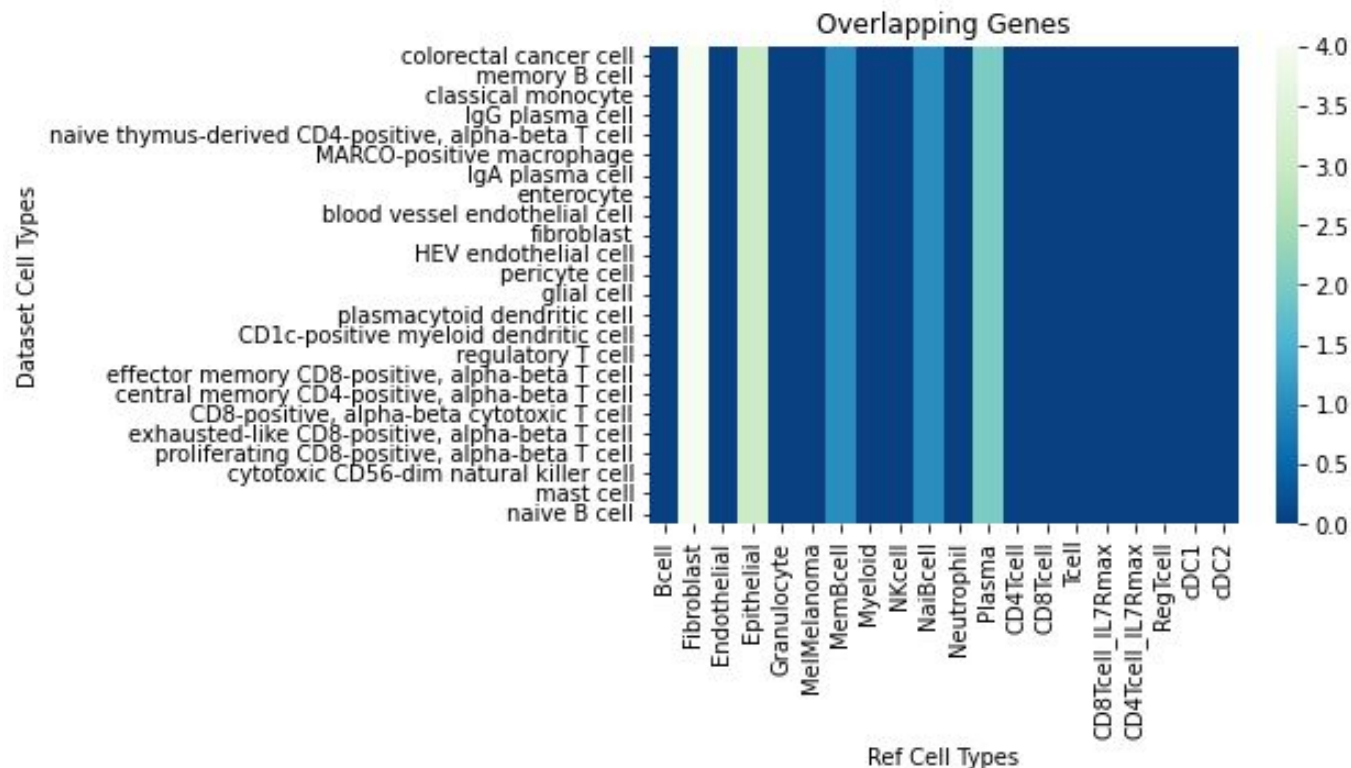
Cell Type Label QC

Algorithm 1: Cell Type Label QC From Reference Marker Genes

input : scRNA-seq AnnData obj. \mathcal{A} , dictionary \mathcal{M} of reference cell type markers from GMT file \mathcal{G}
Hyperparameters $flag_threshold$
output: Dictionary \mathcal{D}_{flag} of Flagged Cell Type Labels
Heatmap \mathcal{H} of Number of Genes

- 1 Initialize dictionary Dictionary \mathcal{D}_{flag} to store QC flagged cell types
- 2 Initialize list $unique_celltypes$ of unique cell types in \mathcal{A}
- 3 **Execute** Differential Gene Expression analysis on \mathcal{A}
- 4 **for** $i \leftarrow 0$ **to** $length(unique_celltypes)$ **do**
- 5 **for** all key value pairs (k, v) in \mathcal{M} **do**
- 6 Initialize list $ref_genes = v$
- 7 Initialize list $genes$ of differentially expressed genes from $unique_celltypes[i]$ in \mathcal{A}
- 8 $S = (genes \cap ref_genes)$
- 9 **if** $length\ S \geq flag_threshold$ **then**
- 10 **Append** S and ref celltype k to $\mathcal{D}_{flag}[unique_celltypes[i]]$
- 11 **end**
- 12 **end**
- 13 **end**
- 14 Create Heatmap, \mathcal{H} with y axis as cell type in \mathcal{D} , x axis as ref cell type in \mathcal{M}
- 15 and cell colors as number of overlapping genes
- 16 **return** $\mathcal{D}_{flag}, \mathcal{H}$

Ex Output - Algorithm 1



Summary

- Differential gene expression analysis as Feature selection
- Association measure variations: Pearson, Spearman, distance correlation, ...
- Clustering algorithms variations: k-means, hierarchical clustering, ...
- Visualization variations: heatmap, Sankey diagram (river flow diagram), ...

Outlook:

- Other association measure variations
- Clustering: mixture modeling,
- Further visualizations
- Workflow adaptability for other datasets

Considerations and limitations

Quality control

- B cell genes expressed in all cells?
- unbalanced datasets (Number of patients, number of cells, ...)

Feature selection

- DGE testing – confounding factors,
- different approaches instead of DGE testing

Biological interpretation

- correlation vs. causation vs. coincidence