

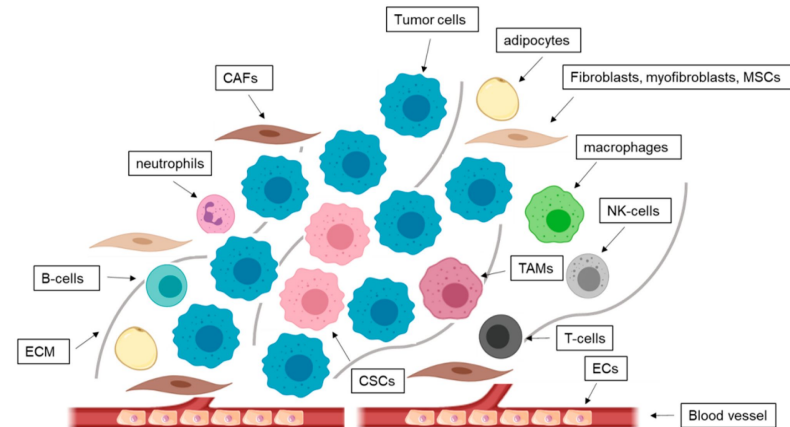
Exploring Human Breast Cancer Atlas scRNA-Seq Data

Team 4

Arthi Hariharan, Seowoo Kim, Taeim Kwon

Biological Concepts behind Questions of Interest

- Tumor microenvironment is heterogenous
 - Immune Cells - T cells, B cells
 - Mesenchymal Stem cells
 - Fibroblasts - Normal, Cancerous
- Breast Cancer Subtypes - HER2+, ER+, TNBC
 - Different Gene Signatures
 - Treatment resistance
 - Gene expression profile - Prognosis



Adapted from Heesch A et al., 2020

Biological Questions:

- Identify the differential expression of genes across breast cancer subtypes.
 - Focus on cancer hallmarks - Metastasis and Inflammation markers (CXCR family) and Angiogenesis markers (VEGF,PDGF family) .
- Show a gene co-expression network in the cells involved with tumor microenvironment, particularly for mesenchymal cells.

Introduction to Dataset

- GSE176078_Wu_et al_2021_BRCA_scRNASeq.tar.gz:
 - A single-cell and spatially resolved atlas of human breast cancers by Wu et al.
 - scRNA-Seq on 26 primary tumors (11 ER+, 5 HER2+ and 10 TNBC)
 - Files
 - Count_matrix_barcodes.tsv
 - Count_matrix_genes.tsv
 - Count_matrix_sparse.mtx
 - Metadata.csv
 - 29,733 features and 100,064 samples

Initial Processing

> project

An object of class Seurat

29733 features across 100064 samples within 1 assay

Active assay: RNA (29733 features, 0 variable features)

1,777 samples were filtered out

> project

An object of class Seurat

29733 features across 98287 samples within 1 assay

Active assay: RNA (29733 features, 2000 variable features)

2 dimensional reductions calculated: pca, umap

Set up Seurat

1. Read in the files
2. Create the Seurat object

Pre-processing

1. Quality control
2. Normalize the data
3. Select the highly variable features
4. Scale the data

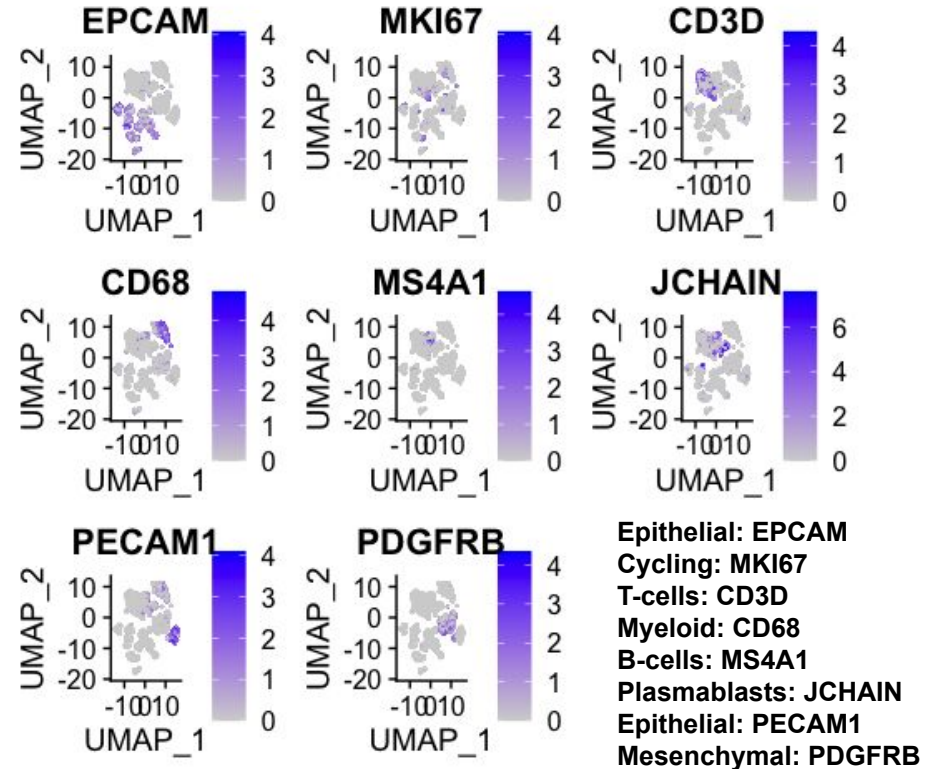
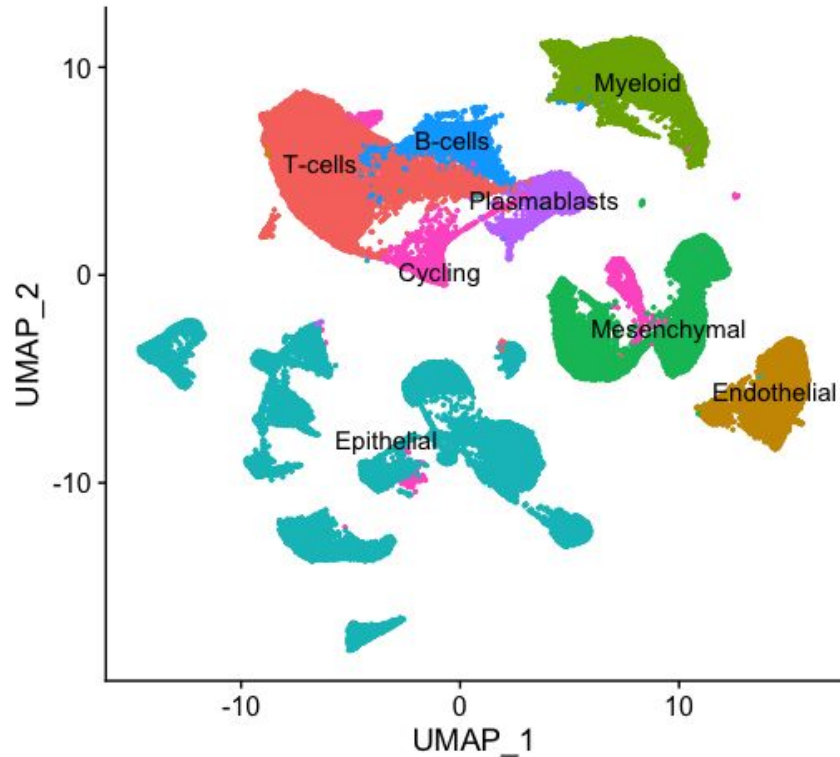
Dimensionality

1. Run PCA
2. Check an 'elbow plot' to choose the dimension of the data
3. Run 'UMAP'

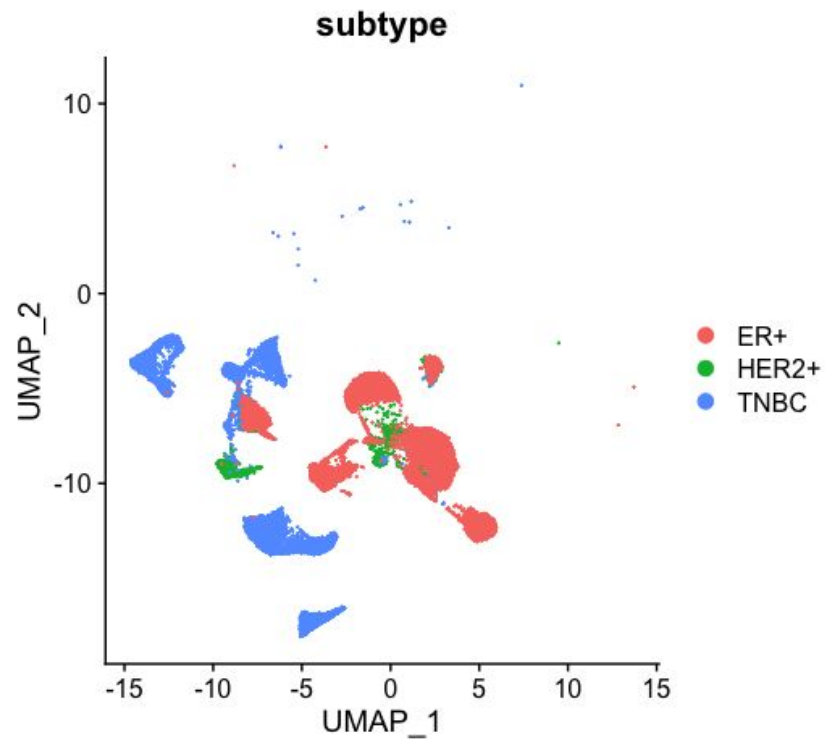
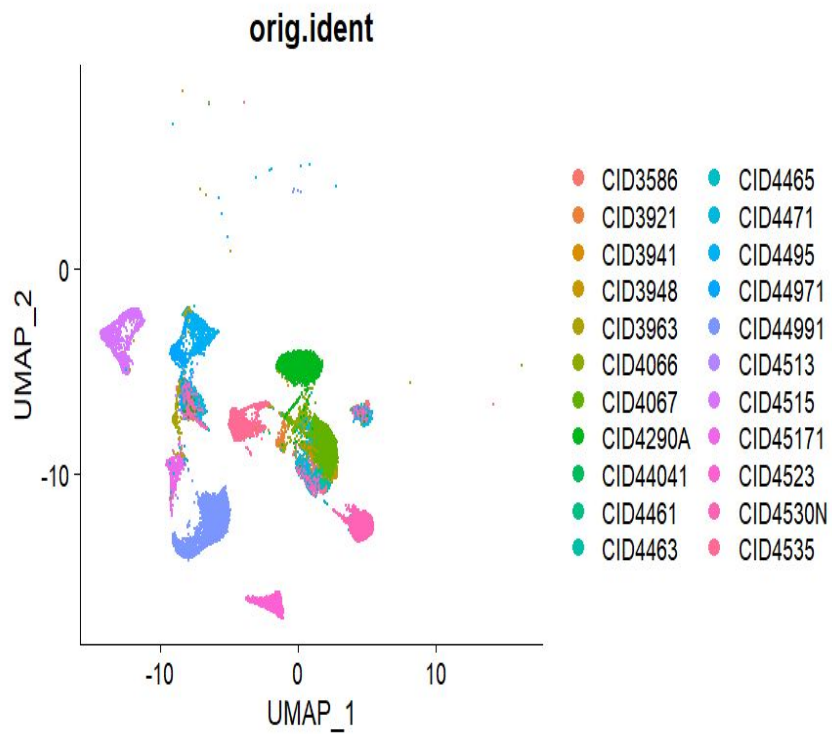
Cell Type Annotation

1. Automatic annotation with the metadata
2. Manual annotation by finding marker genes from each cluster

Reproduced outputs



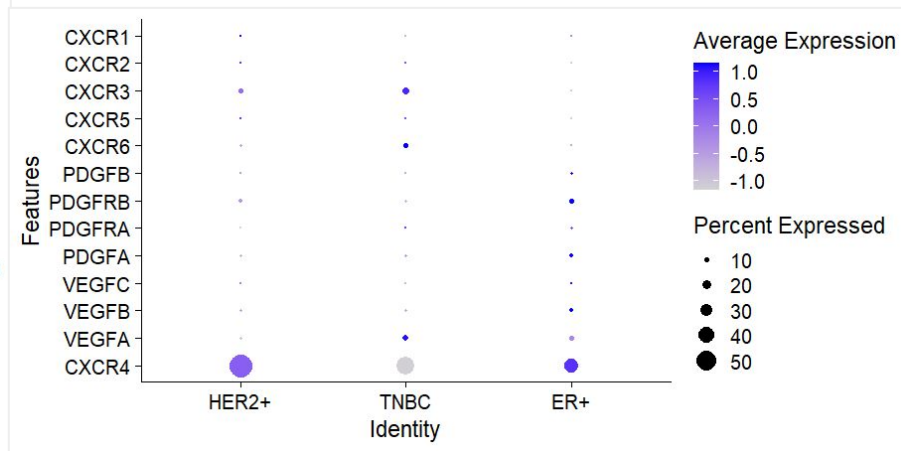
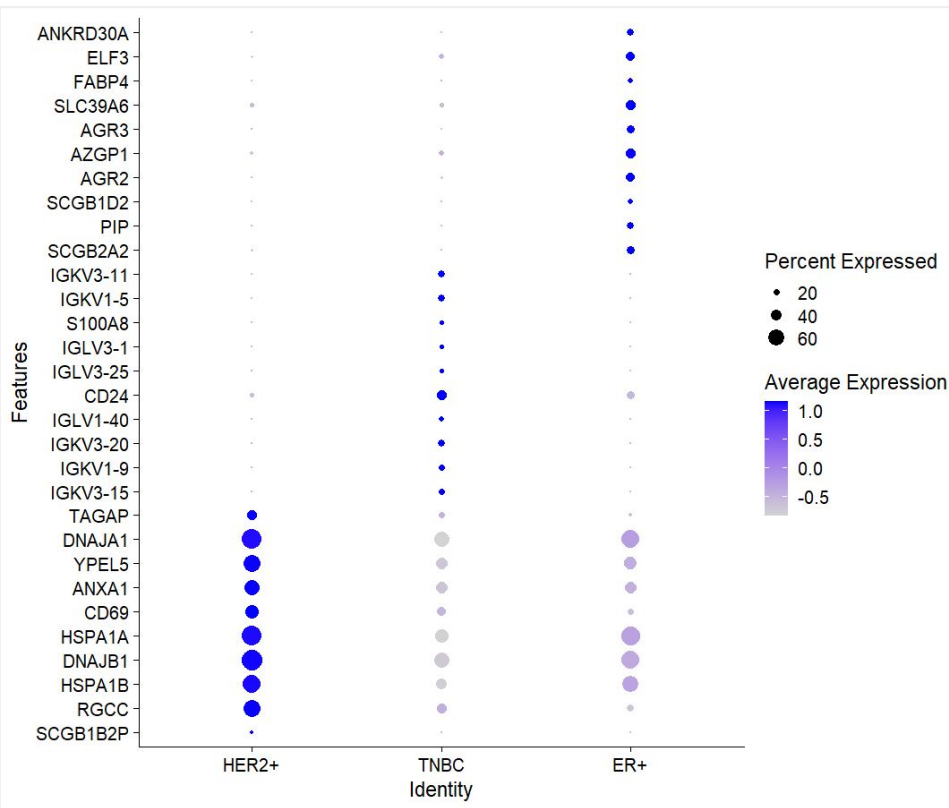
Reproduced outputs



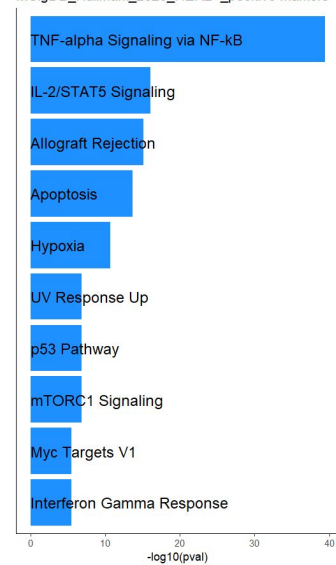
Differential Expression Analysis

- Differential Expression Analysis
 - “FindMarkers” function used to find differentially expressed genes
 - Highly expressed markers in HER2+, ER+ and TNBC
 - CXCR family genes, VEGF, PDGF genes across subtypes
 - EnrichR using MSigDb_hallmarks to find upregulated or downregulated GO terms - “DEEnrichRPlot”

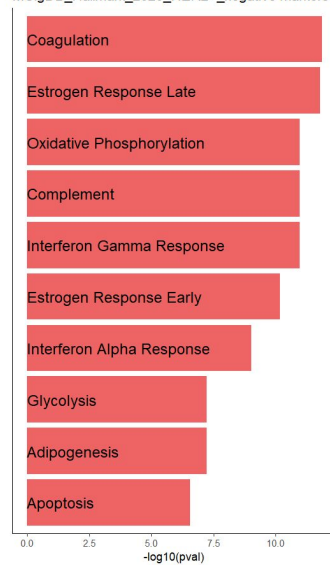
Differential Expression Analysis



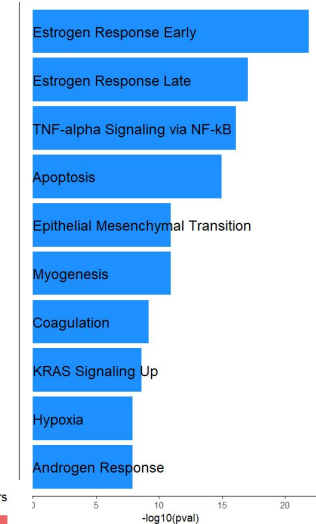
MSigDB_Hallmark_2020_HER2+_positive markers



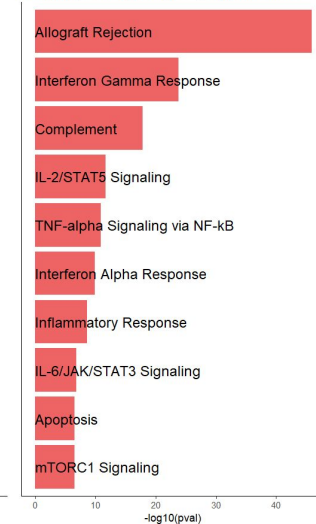
MSigDB_Hallmark_2020_HER2+_negative markers



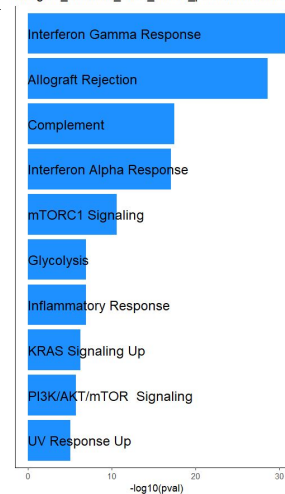
MSigDB_Hallmark_2020_ER+_positive markers



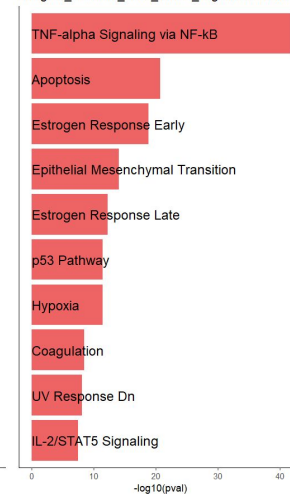
MSigDB_Hallmark_2020_ER+_negative markers



MSigDB_Hallmark_2020_TNBC_positive markers



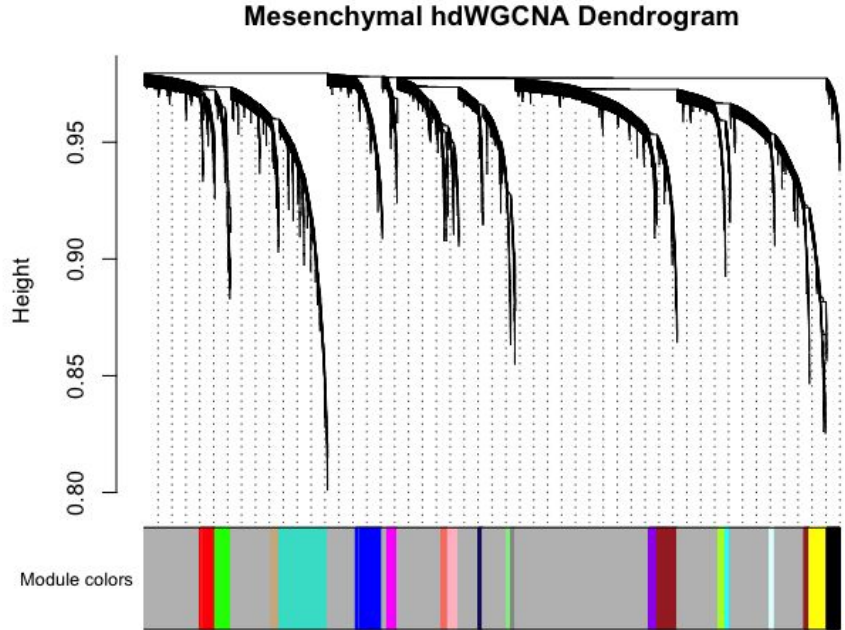
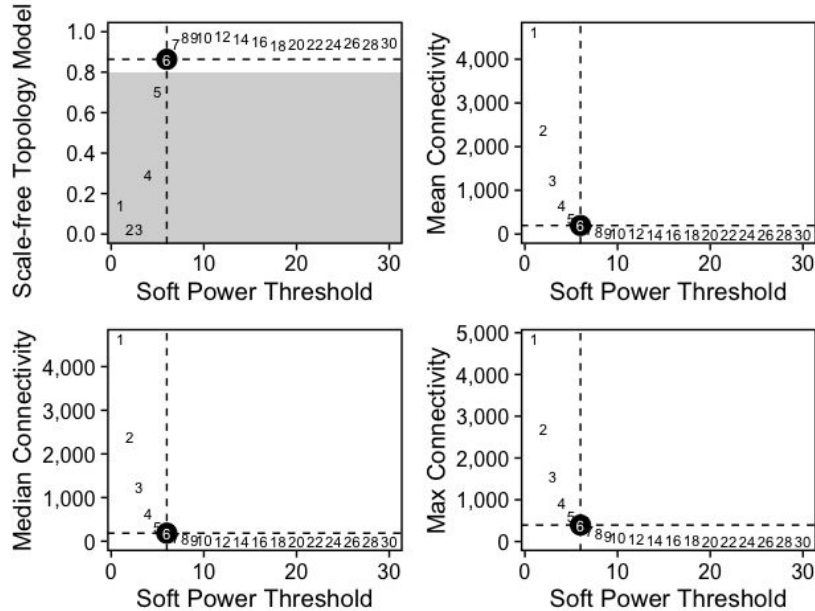
MSigDB_Hallmark_2020_TNBC_negative markers

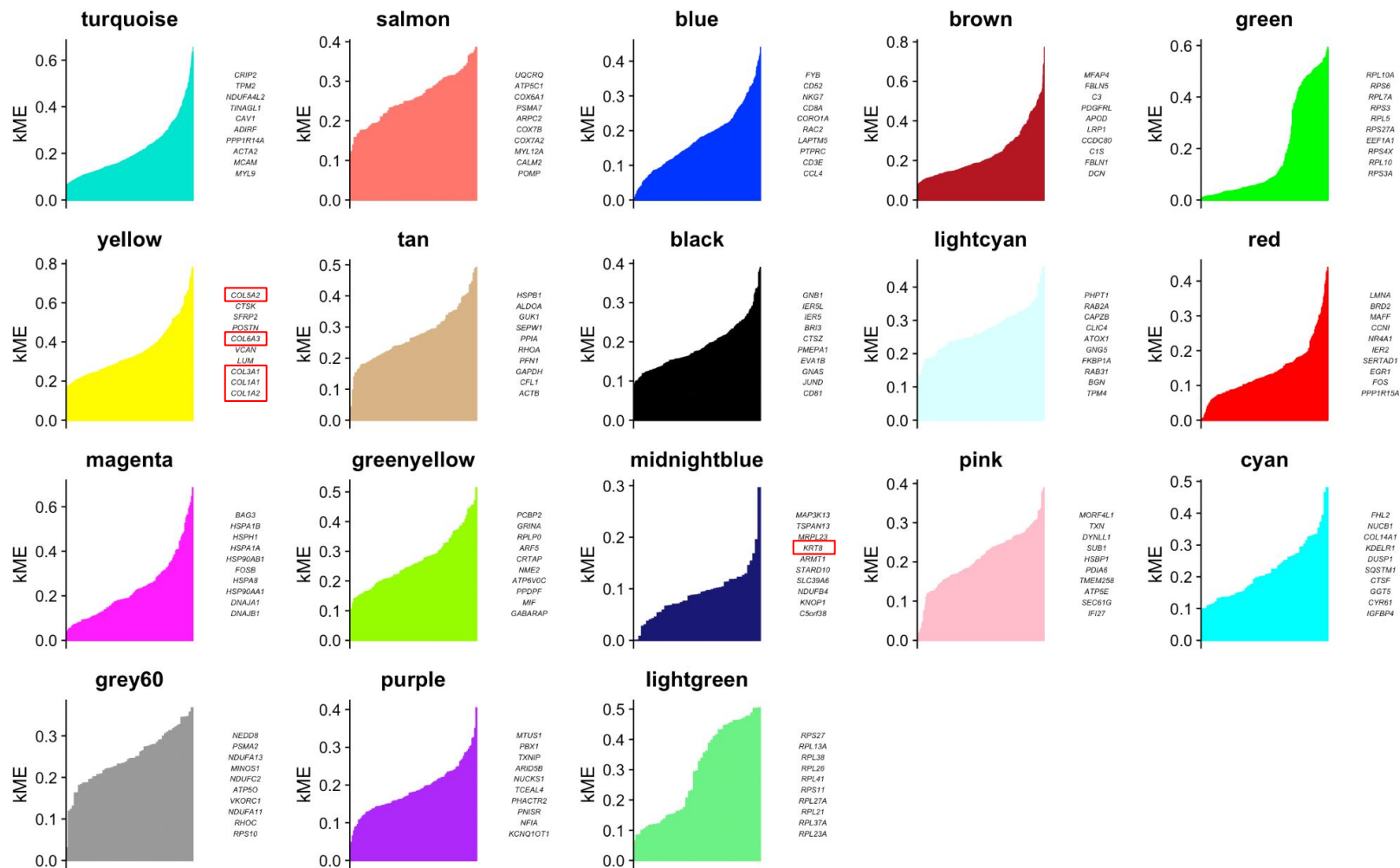


Co-expression Network Analysis

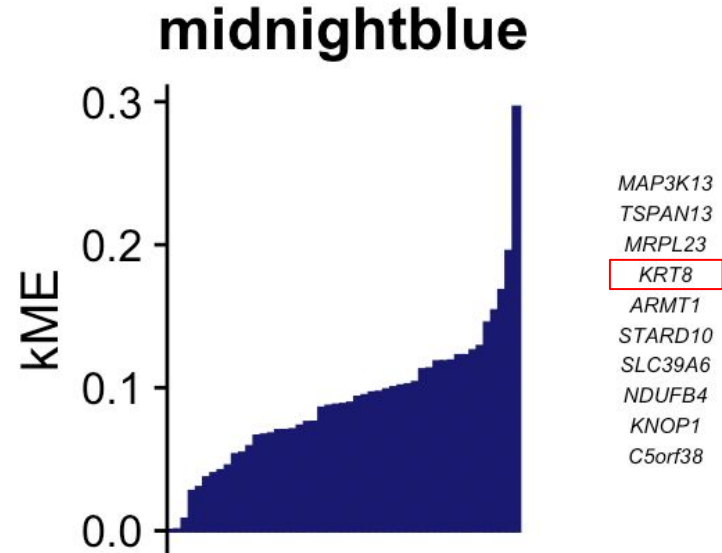
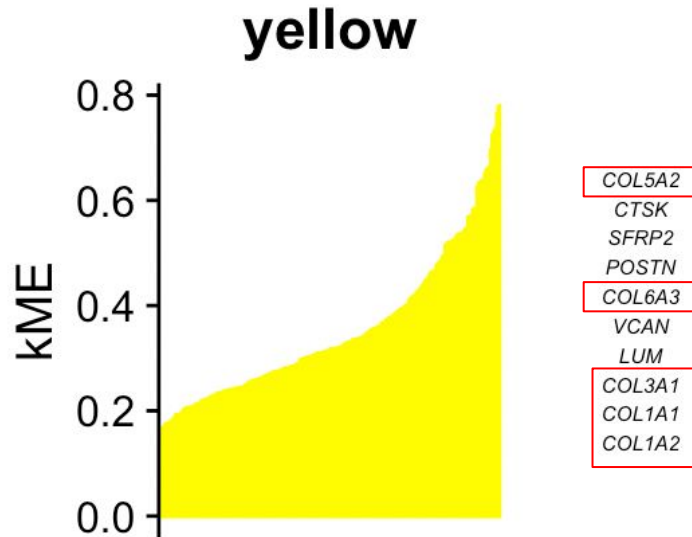
- hdWGCNA
 - Analysis tool for co-expression networks in high dimensional transcriptomics data
 - Co-expression networks
 - Approaches
 - kNN algorithm -> Metacells
 - Hierarchical clustering -> Dendrogram
 - Soft-power threshold
 - Modules
 - PCA -> Module eigengenes
 - Eigengene-based connectivity a.k.a kME
- EnrichR using MSigDb for hallmarks to understand our modules better

Soft-power threshold and Dendrogram



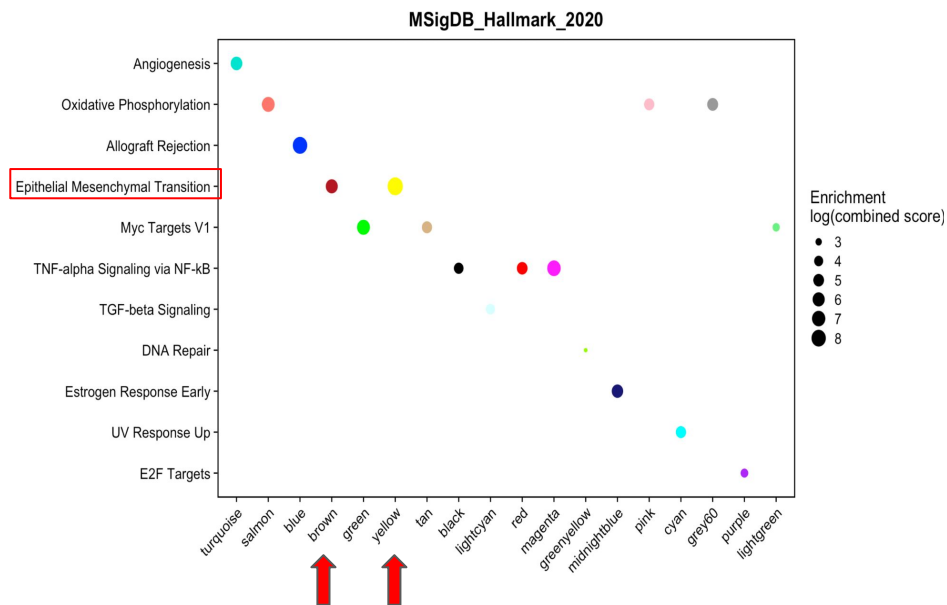


Enlarged plots with ranked genes by each module

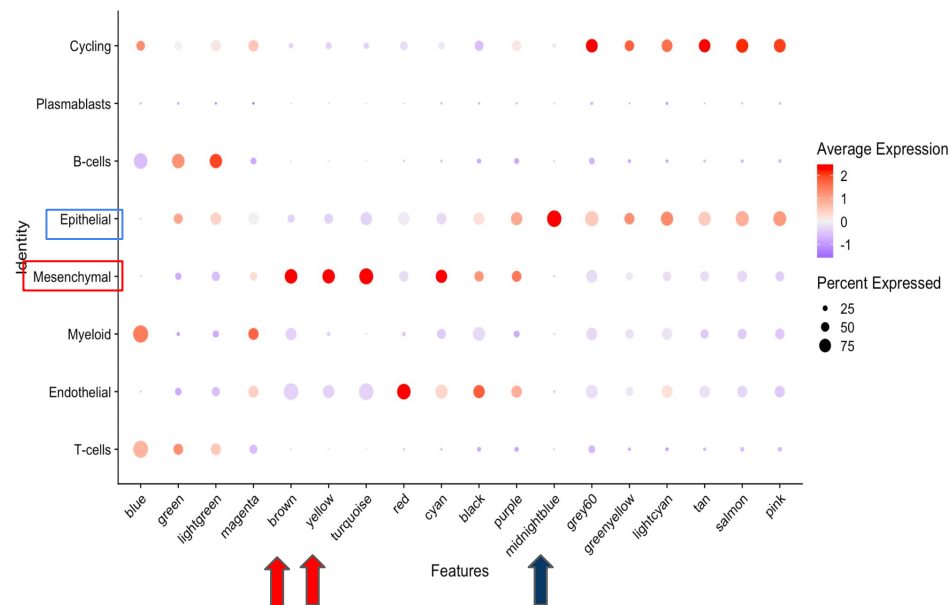


Dot plots for co-expressed gene modules

Across different hallmarks



Across different cell types



Discussion

- **Differential expression**
 - Angiogenesis markers have very low expression across all subtypes
 - Only CXCR4 chemokine is expressed in CXCR gene family
- **Co-expression Network**
 - The module having collagen type genes was highly expressed in the mesenchymal cells
 - The module was enriched in the epithelial mesenchymal transition which is often influenced by tumor microenvironment
- **Contradictory in EnrichR Barplot of differential expressed genes across tumor subtypes**
 - As the dataset was used for spatial transcriptomics, it might cause conflicting results.
 - It would be resolved if we further this analysis taking heterogeneity within the samples into account - divide dataset into specific dimensions or we incorporate the spatial analysis.

References

- Wu, S. Z., Al-Eryani, G., Roden, D. L., Junankar, S., Harvey, K., Andersson, A., Thennavan, A., Wang, C., Torpy, J. R., Bartonicek, N., Wang, T., Larsson, L., Kaczorowski, D., Weisenfeld, N. I., Uytingco, C. R., Chew, J. G., Bent, Z. W., Chan, C. L., Gnanasambandapillai, V., Dutertre, C. A., ... Swarbrick, A. (2021). A single-cell and spatially resolved atlas of human breast cancers. *Nature genetics*, 53(9), 1334–1347. <https://doi.org/10.1038/s41588-021-00911-1>
- Morabito, S., Reese, F., Rahimzadeh, N., Miyoshi, E., & Swarup, V. (2022). High dimensional co-expression networks enable discovery of transcriptomic drivers in complex biological systems. *bioRxiv*, 2022-09.
- Morabito, S., Miyoshi, E., Michael, N., Shahin, S., Martini, A. C., Head, E., ... & Swarup, V. (2021). Single-nucleus chromatin accessibility and transcriptomic characterization of Alzheimer's disease. *Nature genetics*, 53(8), 1143-1155.

Supplementary Plots

