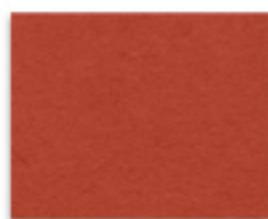


CLUSTERING IN RNA-Seq DATA USING GRADE of MEMBERSHIP MODELS

Dey KK, Hsiao CJ, Stephens M (2017) Visualizing the structure of RNA-seq expression data using grade of membership models. PLOS Genetics 13(3): e1006599. <https://doi.org/10.1371/journal.pgen.1006599>



cell type 1

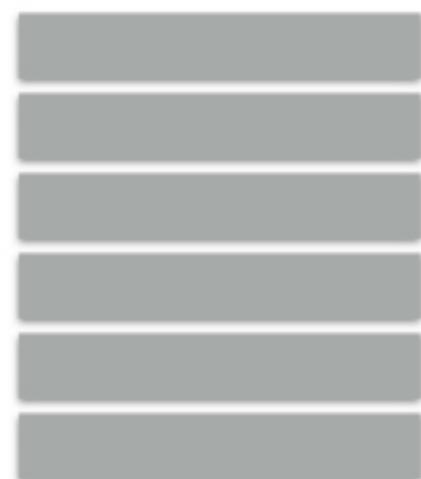
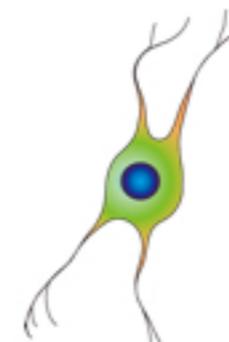
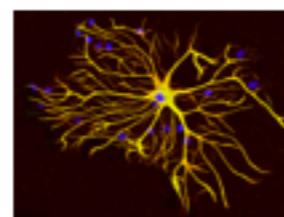
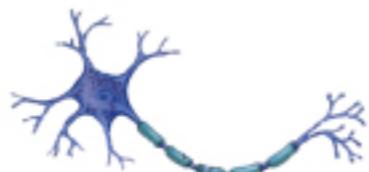


cell type 2

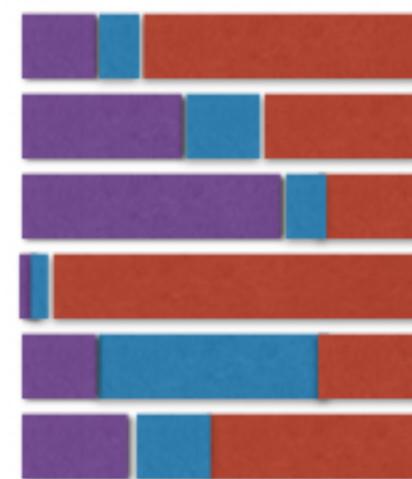


cell type 3

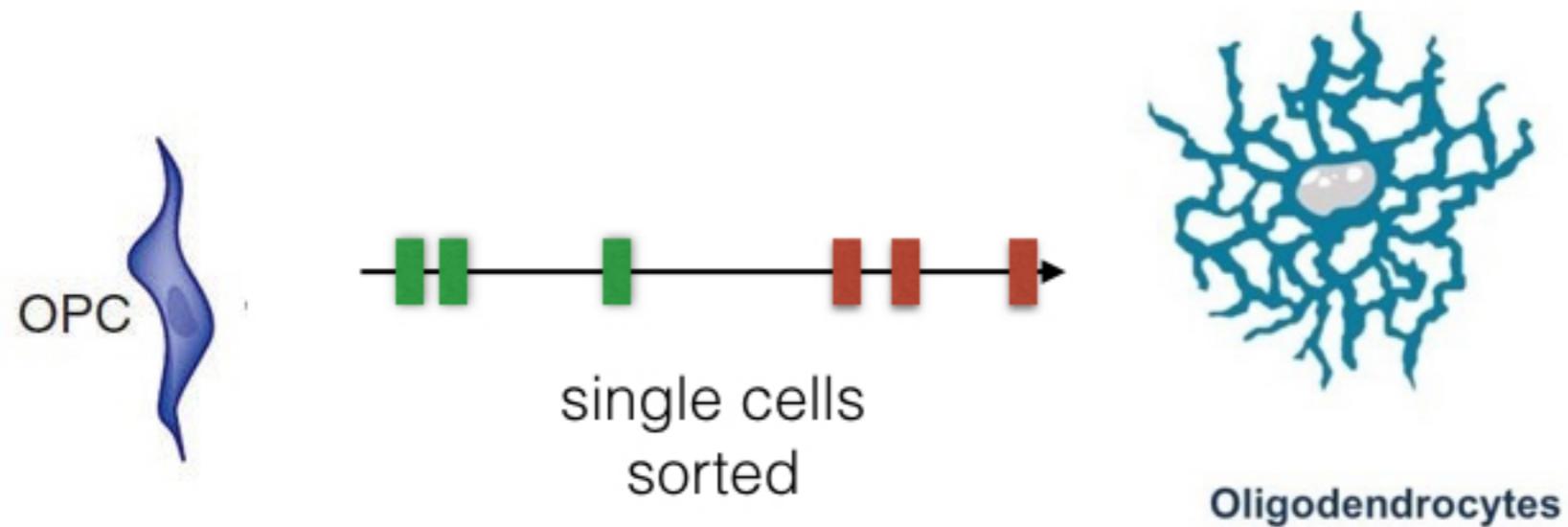
example:
brain



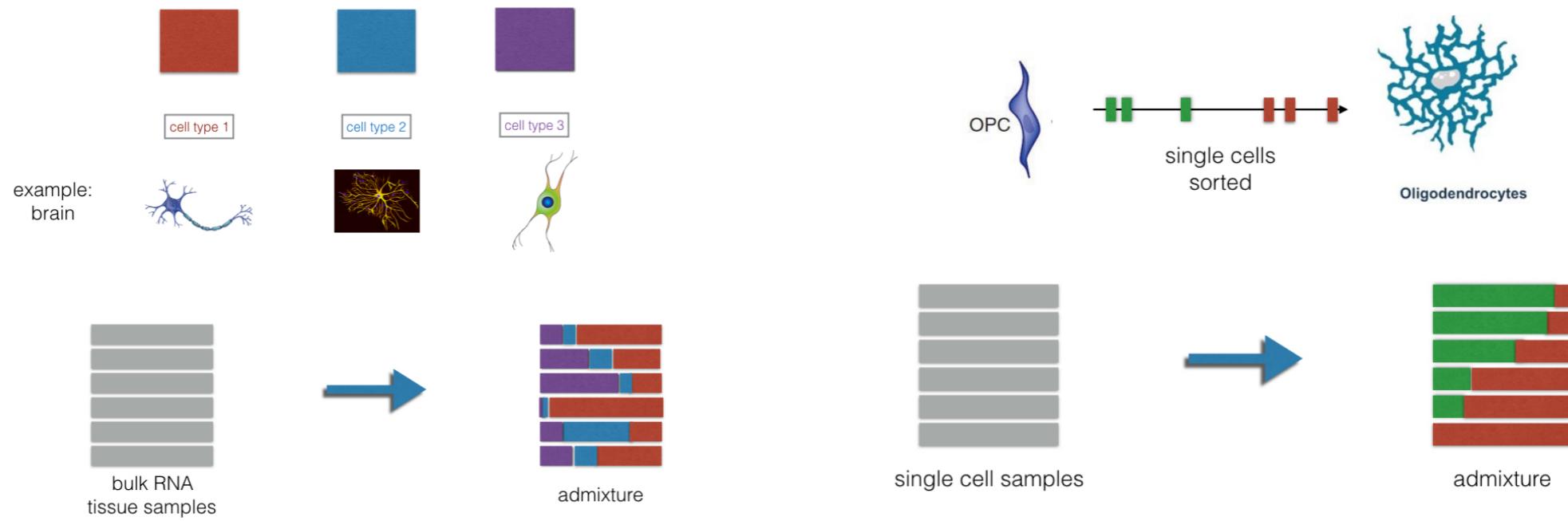
bulk RNA
tissue samples



admixture



Clustering of RNA-seq data

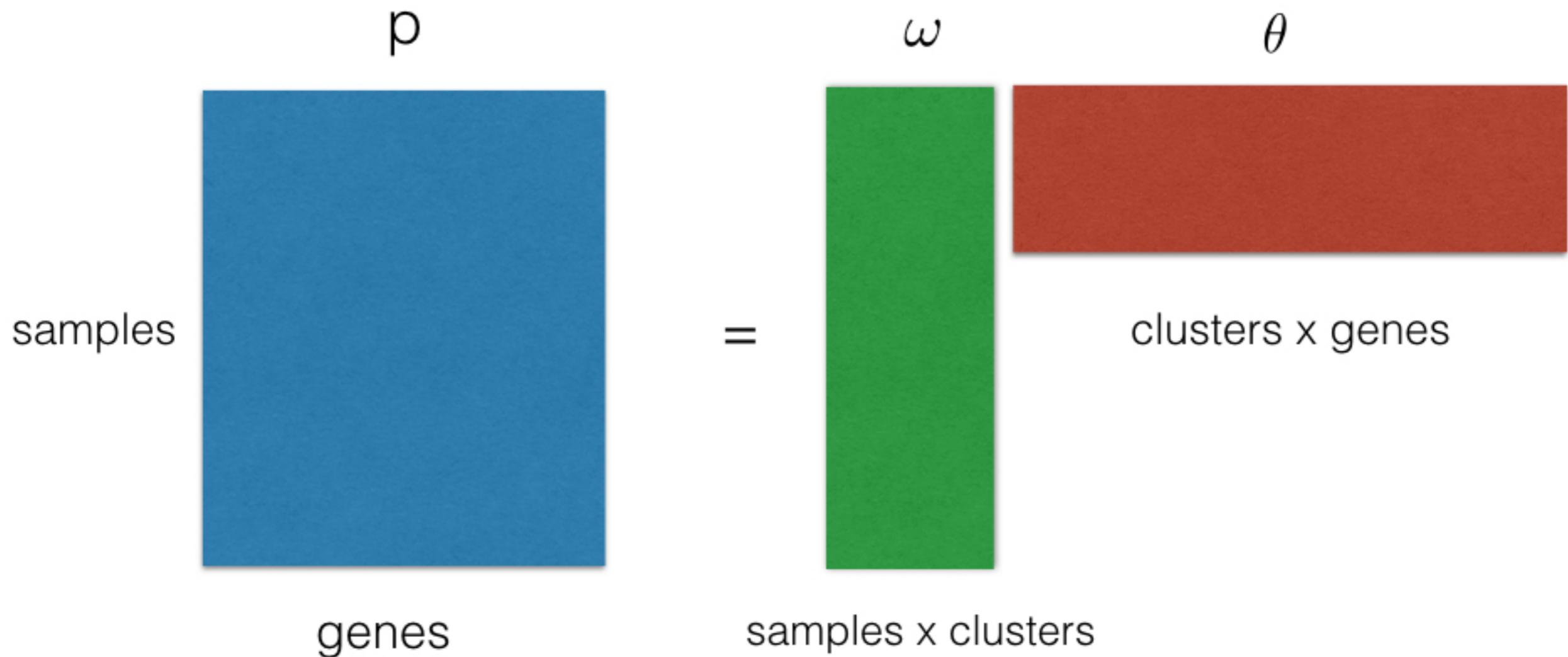


We suggest a model based clustering algorithm in lines of the “admixture” model in population genetics where each sample may have memberships in multiple clusters, which ideally we want to mimic the cell types or cell stages.

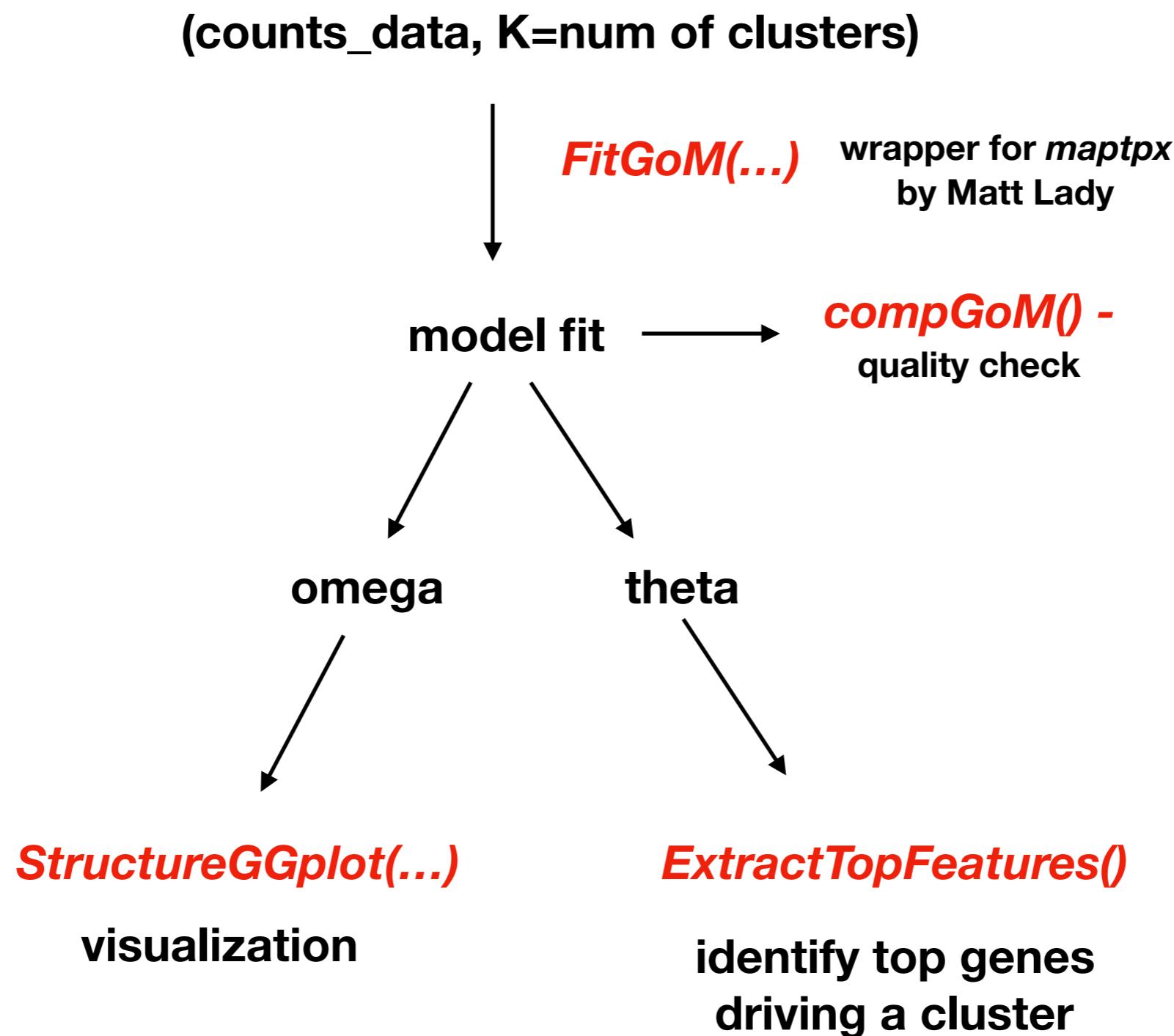
Grade of Membership (GoM) Model

$$(c_{n1}, c_{n2}, \dots, c_{nG}) \sim Mult(c_{n+}, p_{n1}, p_{n2}, \dots, p_{nG})$$

$$p_{ng} = \sum_{k=1}^K \omega_{nk} \theta_{kg} \quad \sum_{g=1}^G \theta_{kg} = 1 \quad \forall k \quad \sum_{k=1}^K \omega_{nk} = 1 \quad \forall n$$



Model workflow (R package CountClust)



How *FitGoM* works?

We assume Dirichlet priors on the parameters ω and θ
both cases with equal mean for all the components

The initial choices of omega and theta are randomly chosen from the prior

An EM type algorithm is run to update the omega and theta until convergence in log posterior probability.

log Bayes factor corresponding to the model and log BIC are reported to check for model fit and compare between model fits (*compGoM*)

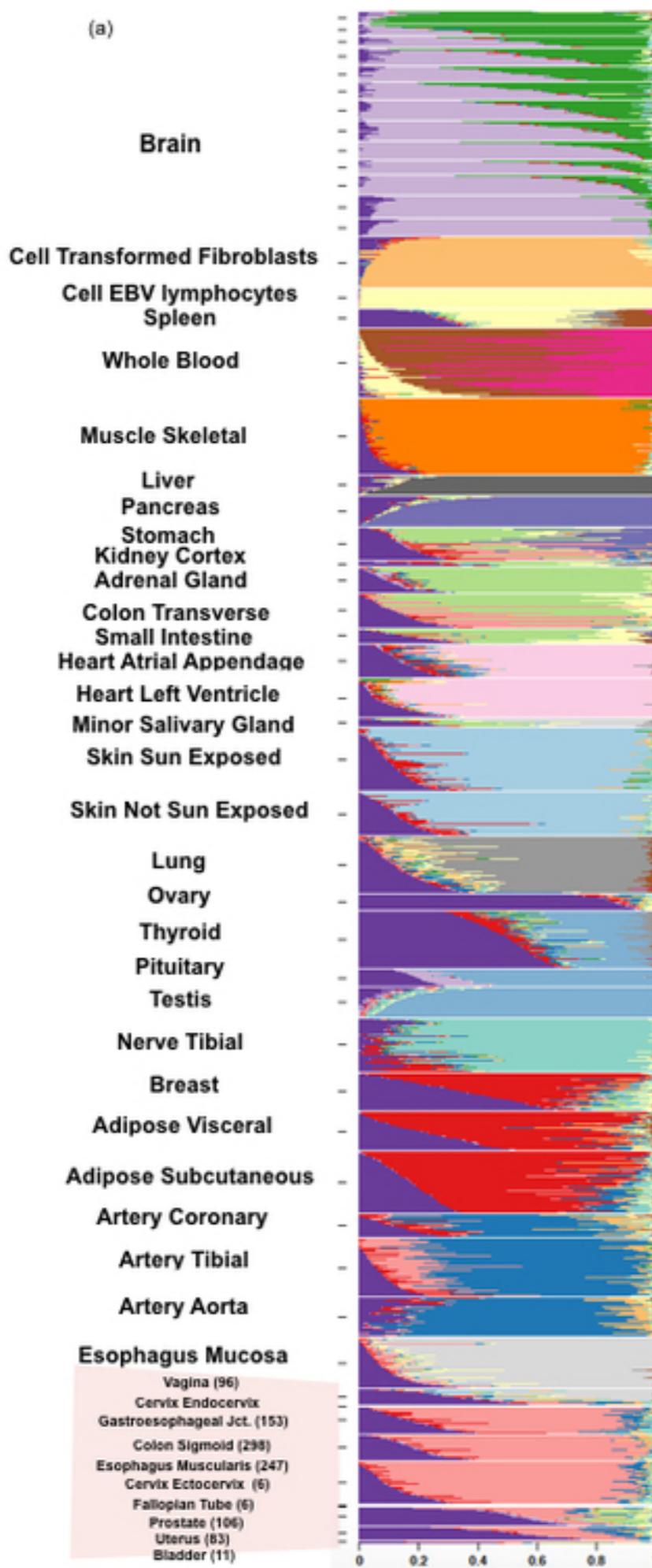
How *ExtractTopFeatures* works?

$$D^g(k) := \min_{l:l \neq k} \left[\theta_{kg} \log \frac{\theta_{kg}}{\theta_{lg}} + \theta_{lg} - \theta_{kg} \right]$$

We compute the above for each gene g and cluster k

For each k , we identify the genes with highest value of the above measure, and which also have the highest value of θ of all the clusters.

(a)



ExtractTopFeatures

Brown cluster : imp genes

CSF3R, MMP25, IL1R2, SELL

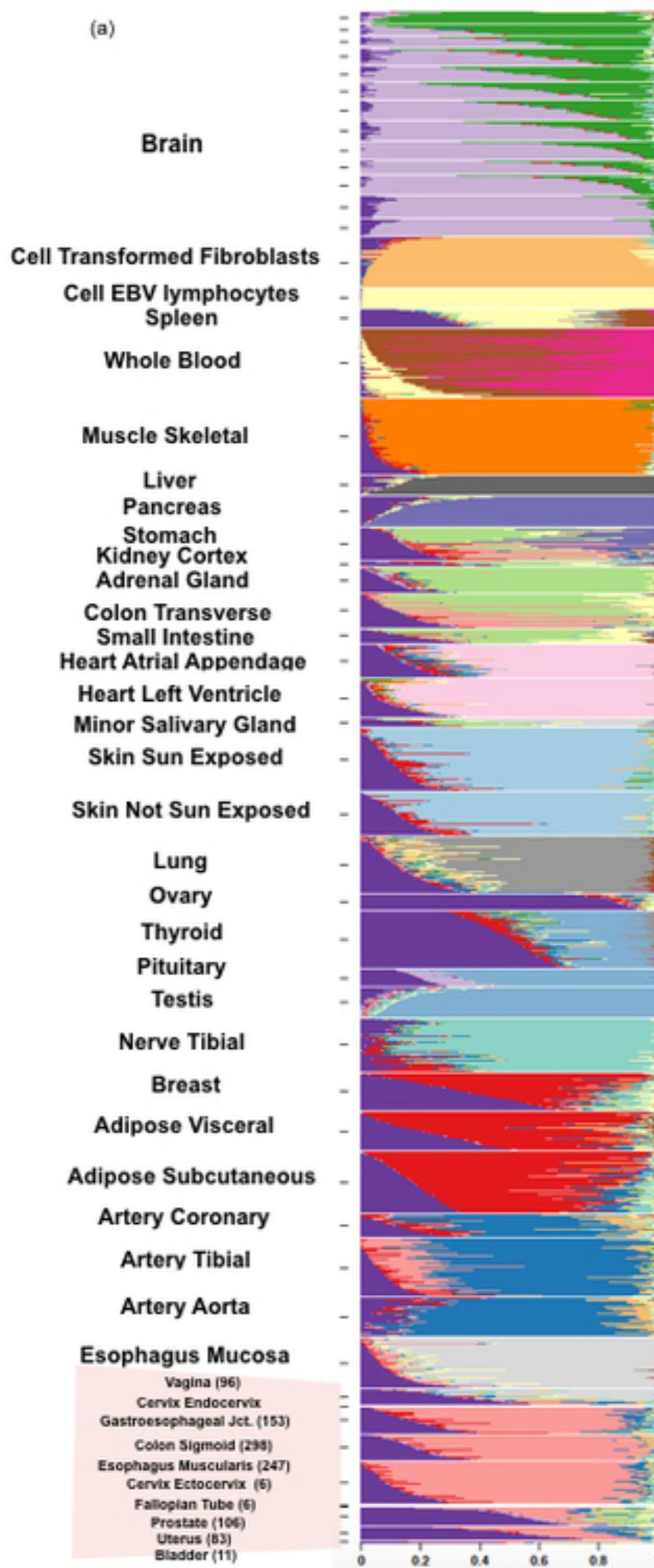
Enriched GO: immune response, defense response, cell periphery, regulation of immune response

Pink cluster: imp genes

HBB, HBA1, HBA2, HBD

Enriched GO: hemoglobin complex, heme binding, gas transport

(a)



ExtractTopFeatures

Light blue cluster : imp genes
PRSS1, CPA1, PNLIP, CELA3A

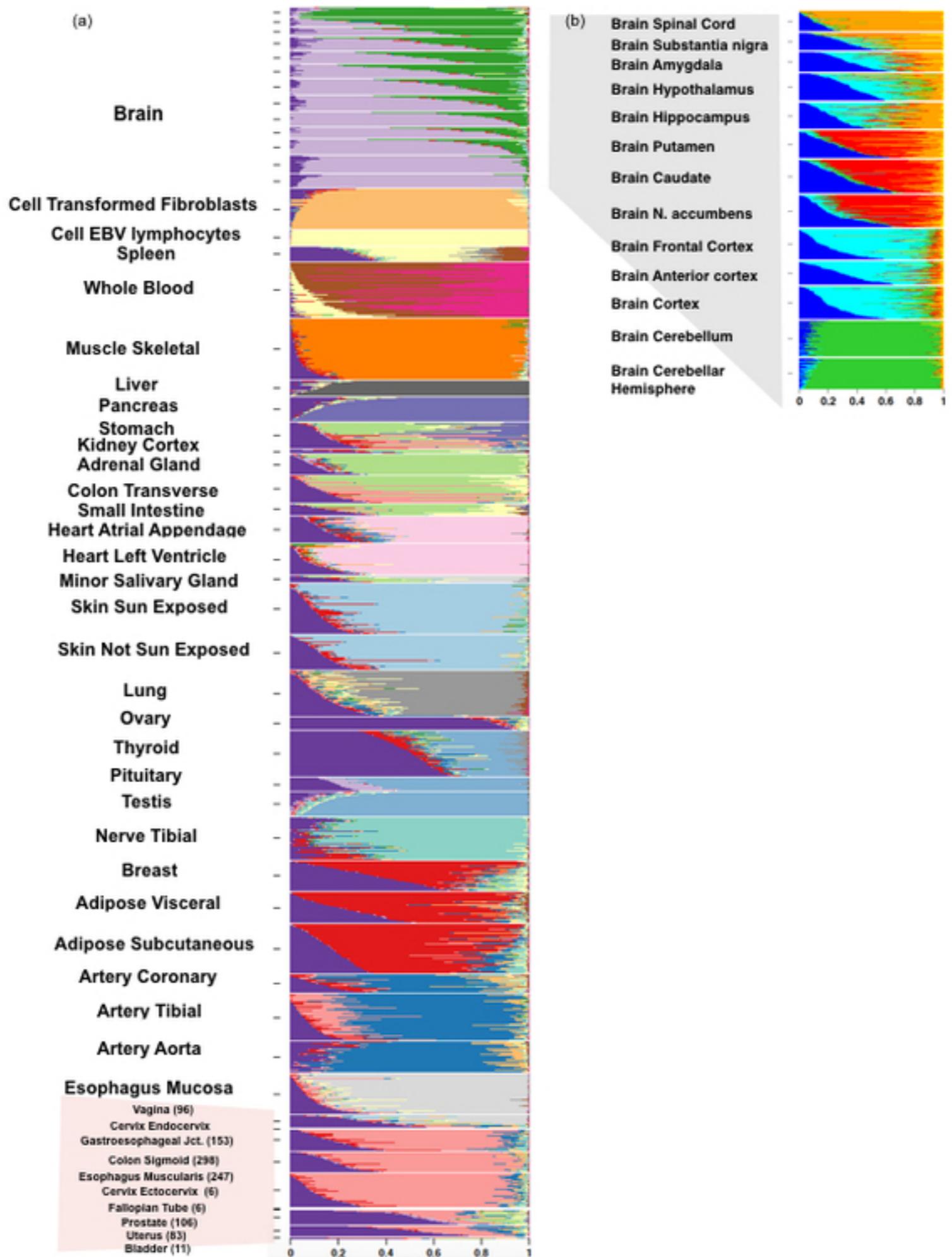
Enriched GO: digestion, serine type endopeptidase activity, proteolysis, hydrolase activity

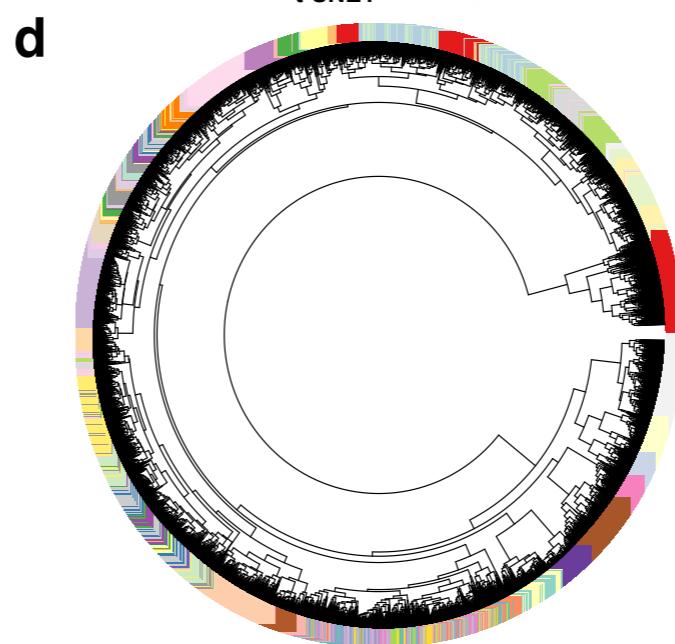
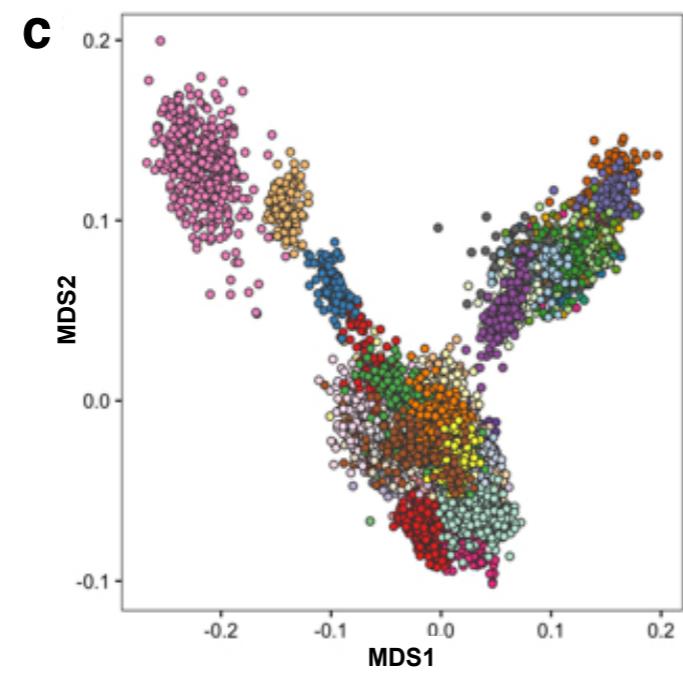
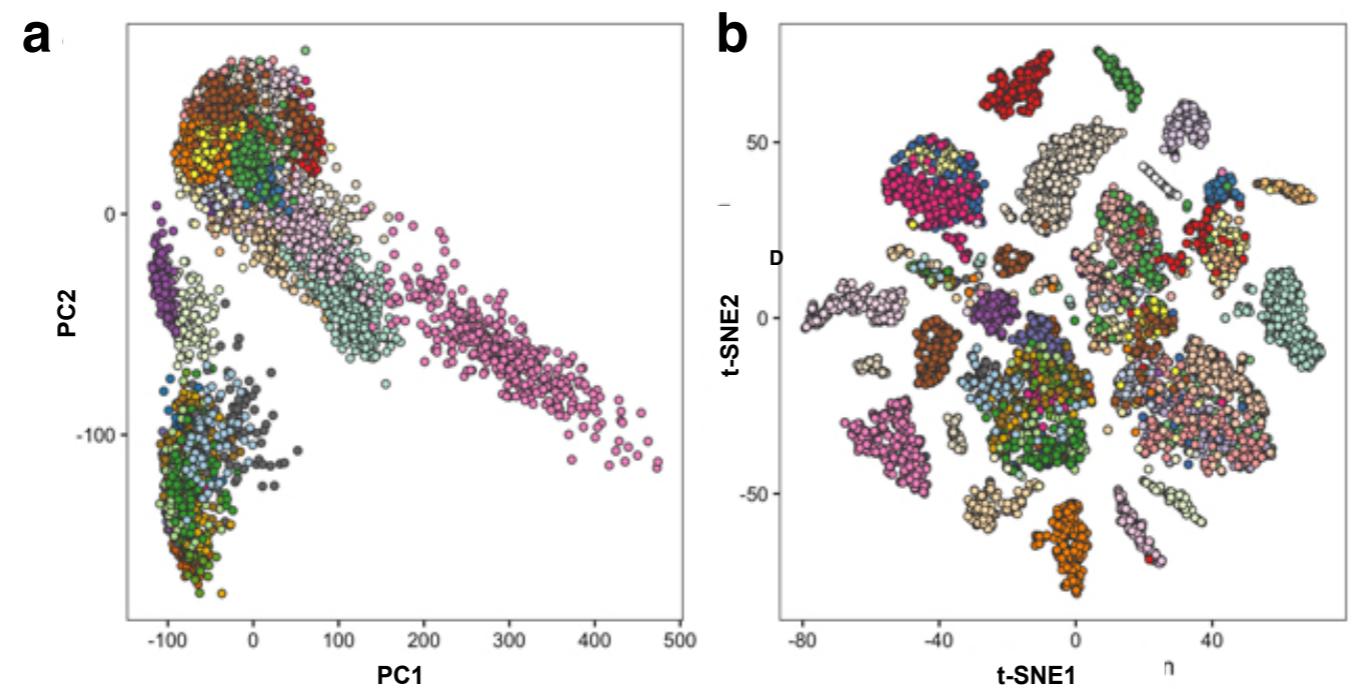
Light blue cluster : imp genes
KRT10, KRT1, KRT2, LOR

Enriched GO: epidermis development, molting cycle, hair cycle

Red cluster: imp genes
FABP4, PLIN1, FASN, GPX3

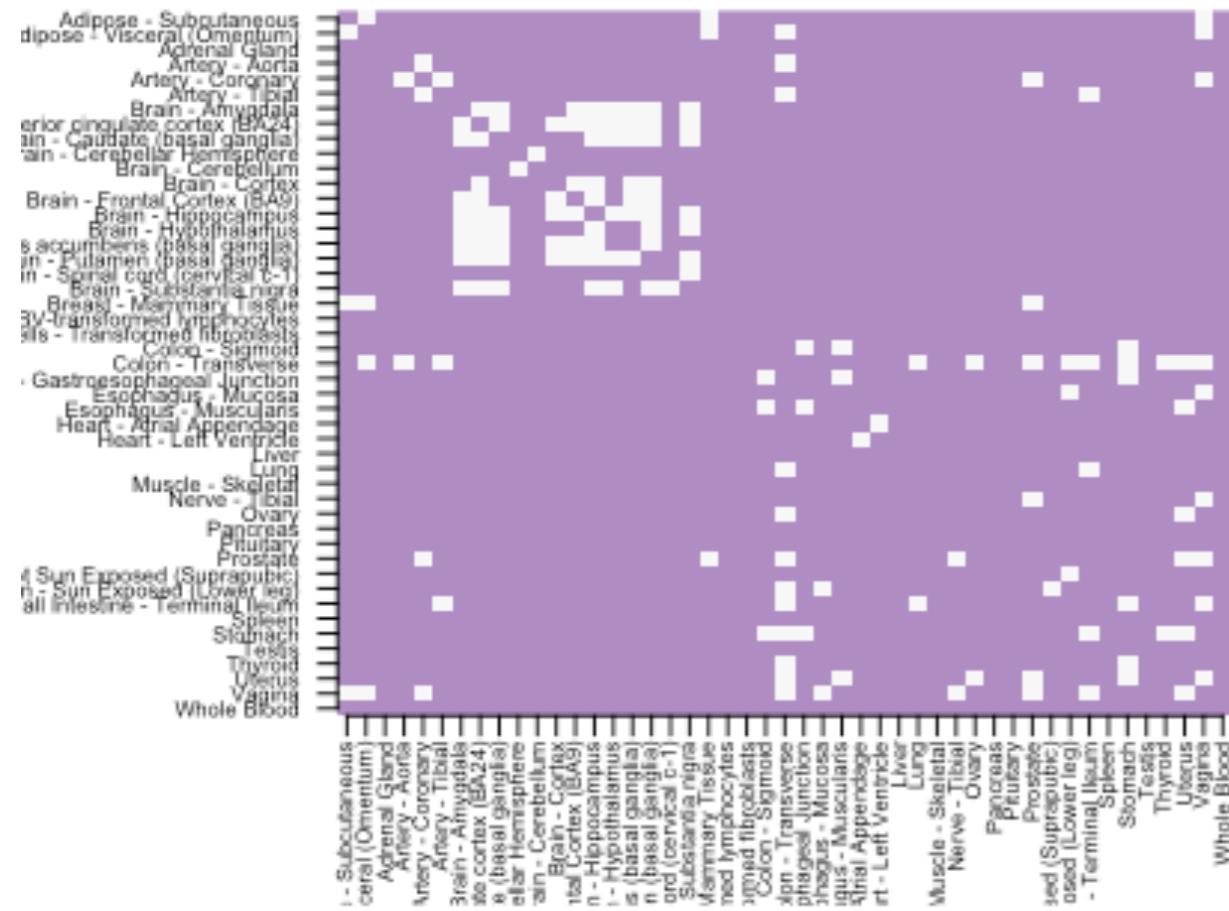
Enriched GO: cellular lipid metabolism, acylglycerol metabolism, angiogenesis regulation



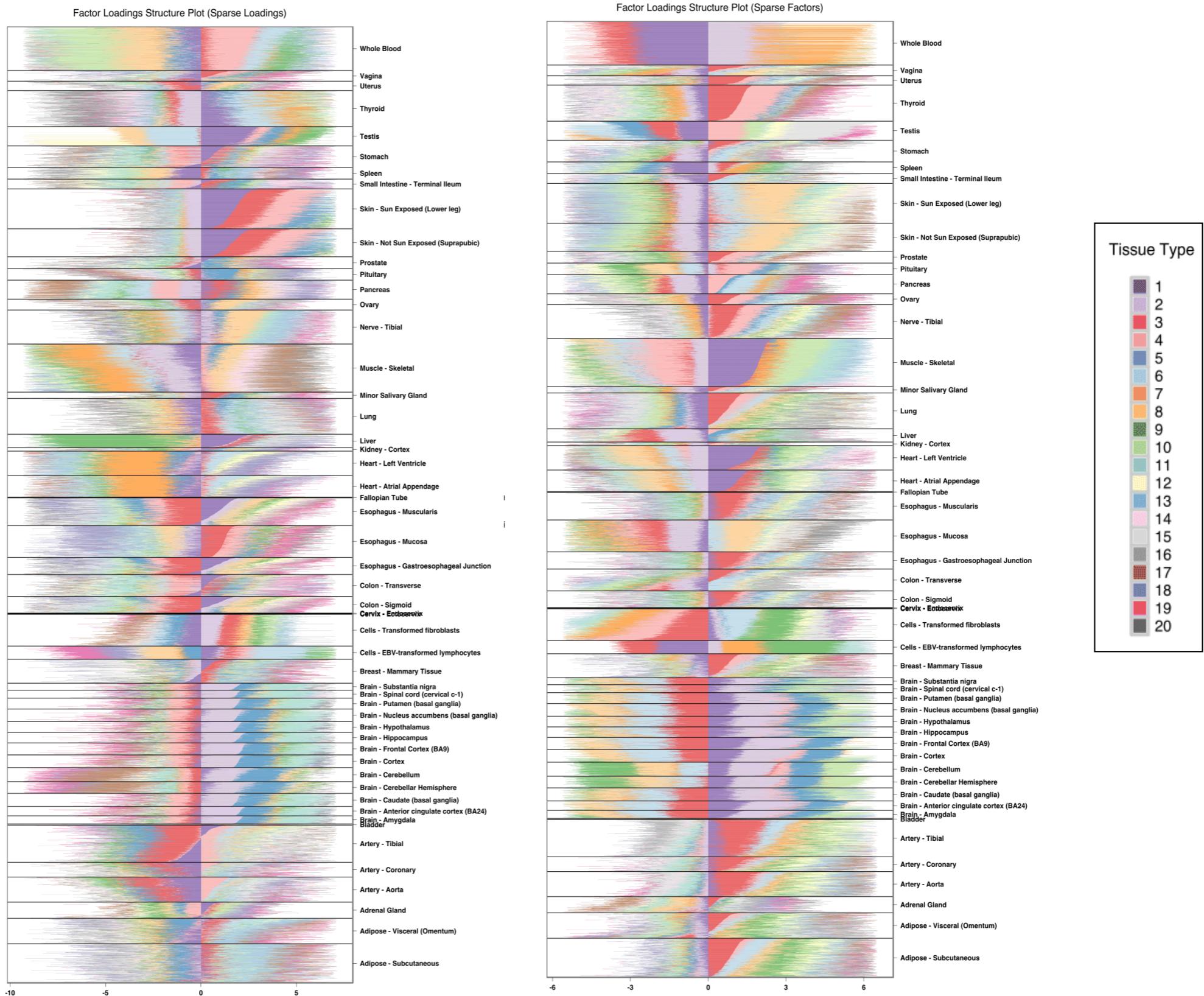


Adipose - Subcutaneous	Brain - Hippocampus	Colon - Transverse	Ovary
Adipose - Visceral (Omentum)	Brain - Hypothalamus	Esophagus - Mucosa	Pancreas
Adrenal Gland	Brain - Spinal cord (cervical c-1)	Esophagus - Muscularis	Pituitary
Artery - Aorta	Brain - Substantia nigra	Esophagus - Gastroesophageal Jn.	Prostate
Artery - Coronary	Brain - Caudate	Fallopian Tube	Skin - Sun Exposed (Lower leg)
Artery - Tibial	Brain - N. accumbens	Heart - Atrial Appendage	Skin - Unexposed (Suprapubic)
Bladder	Brain - Putamen	Heart - Left Ventricle	Small Intestine - Terminal Ileum
Brain - Amygdala	Breast - Mammary Tissue	Kidney - Cortex	Spleen
Brain - Anterior cortex (BA24)	Cells - Transformed fibroblasts	Liver	Stomach
Brain - Cerebellar Hemisphere	Cells - EBV-lymphocytes	Lung	Testis
Brain - Cerebellum	Cervix - Ectocervix	Minor Salivary Gland	Thyroid
Brain - Cortex	Cervix - Endocervix	Muscle - Skeletal	Uterus
Brain - Frontal Cortex (BA9)	Colon - Sigmoid	Nerve - Tibial	Vagina
			Whole Blood

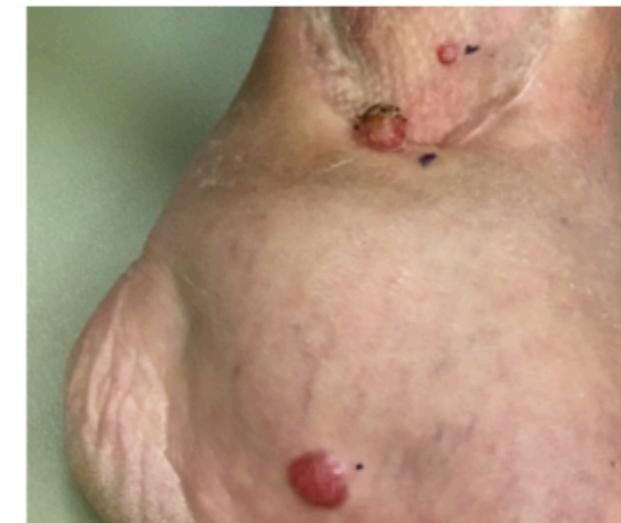
Comparison - GoM vs Hierarchical Method



Sparse Factor Analysis : Engelhardt and Stephens (2010)

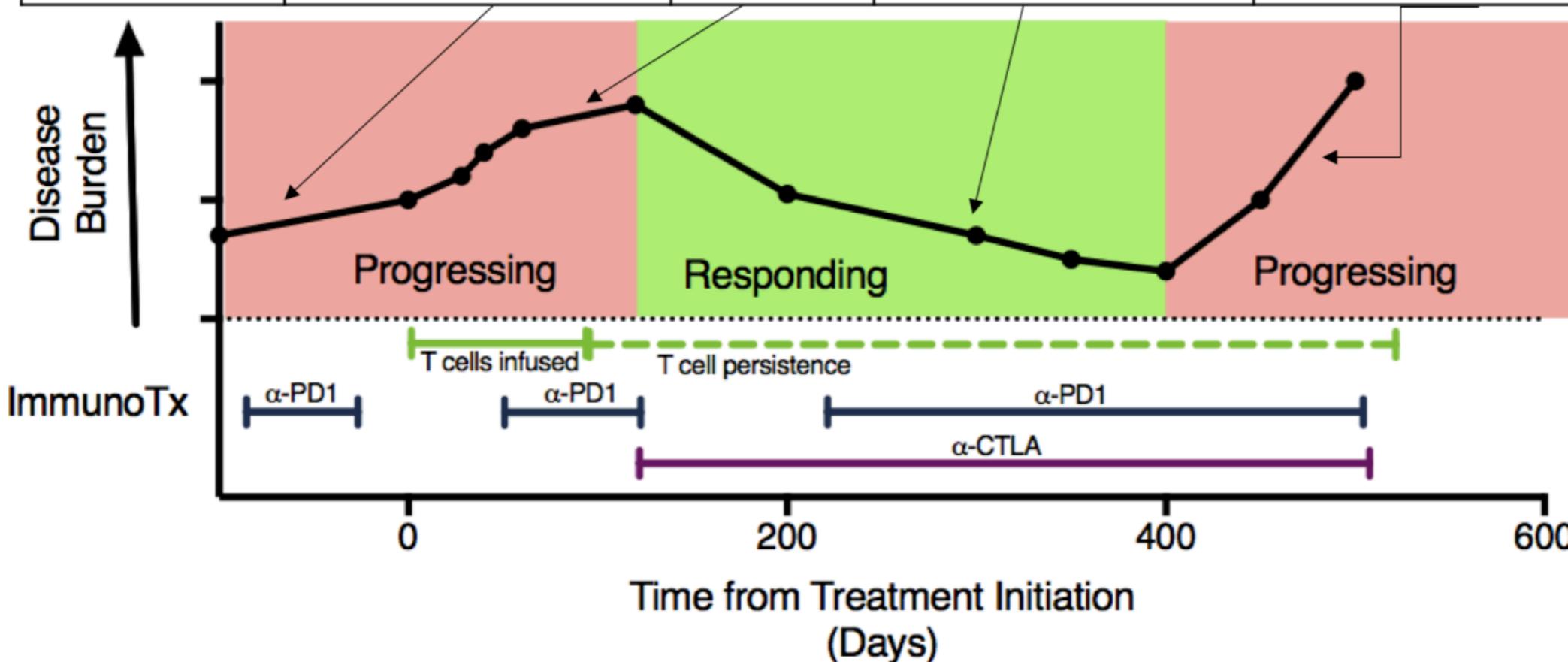


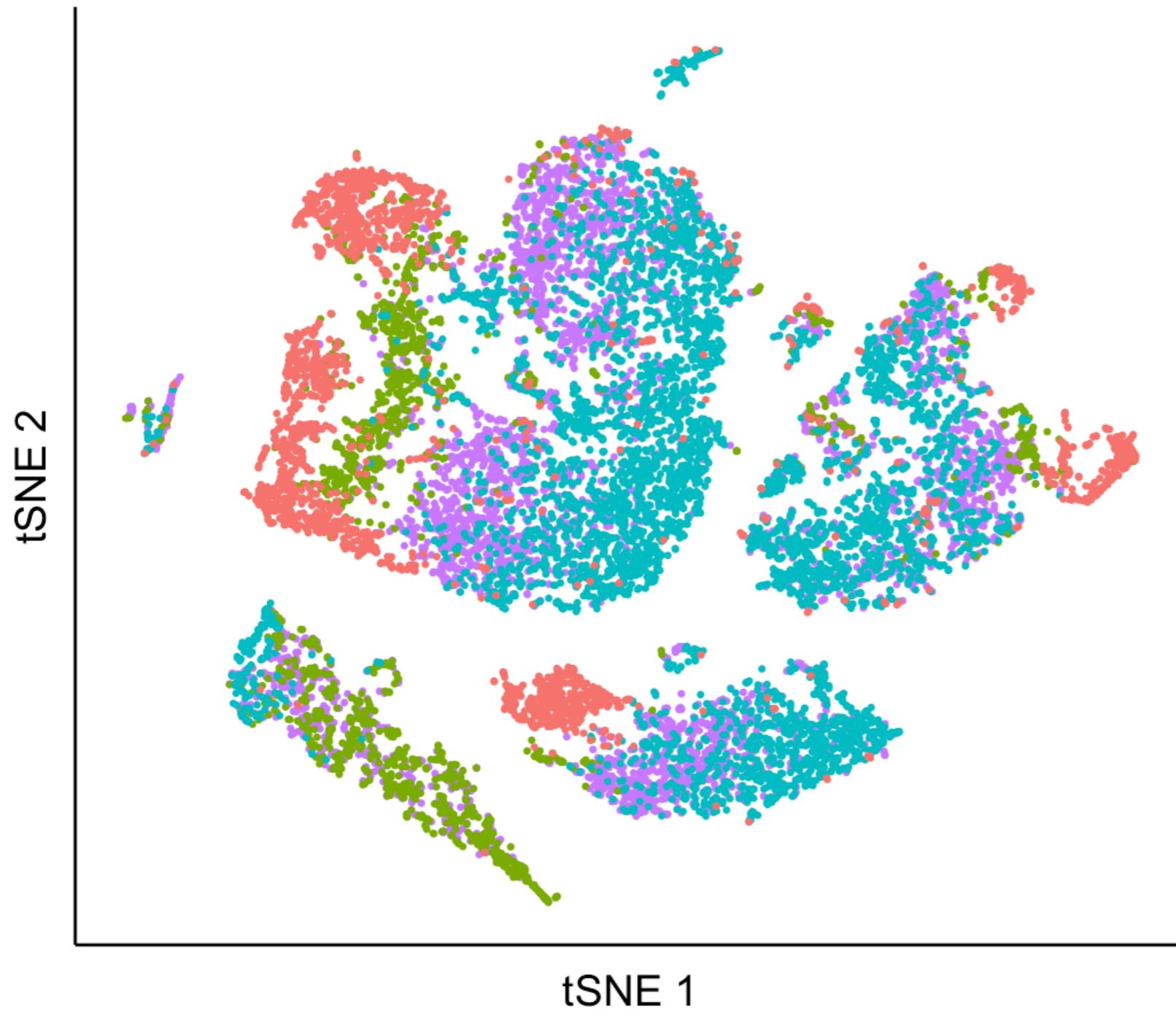
Merkel Cell Carcinoma Study



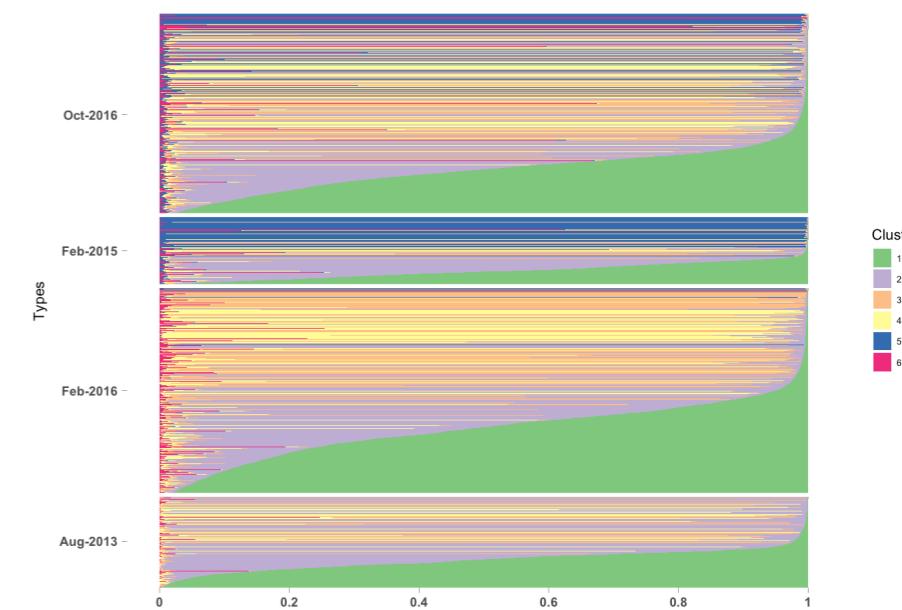
Aug 2013 Feb 2015 Feb 2016 Oct 2016

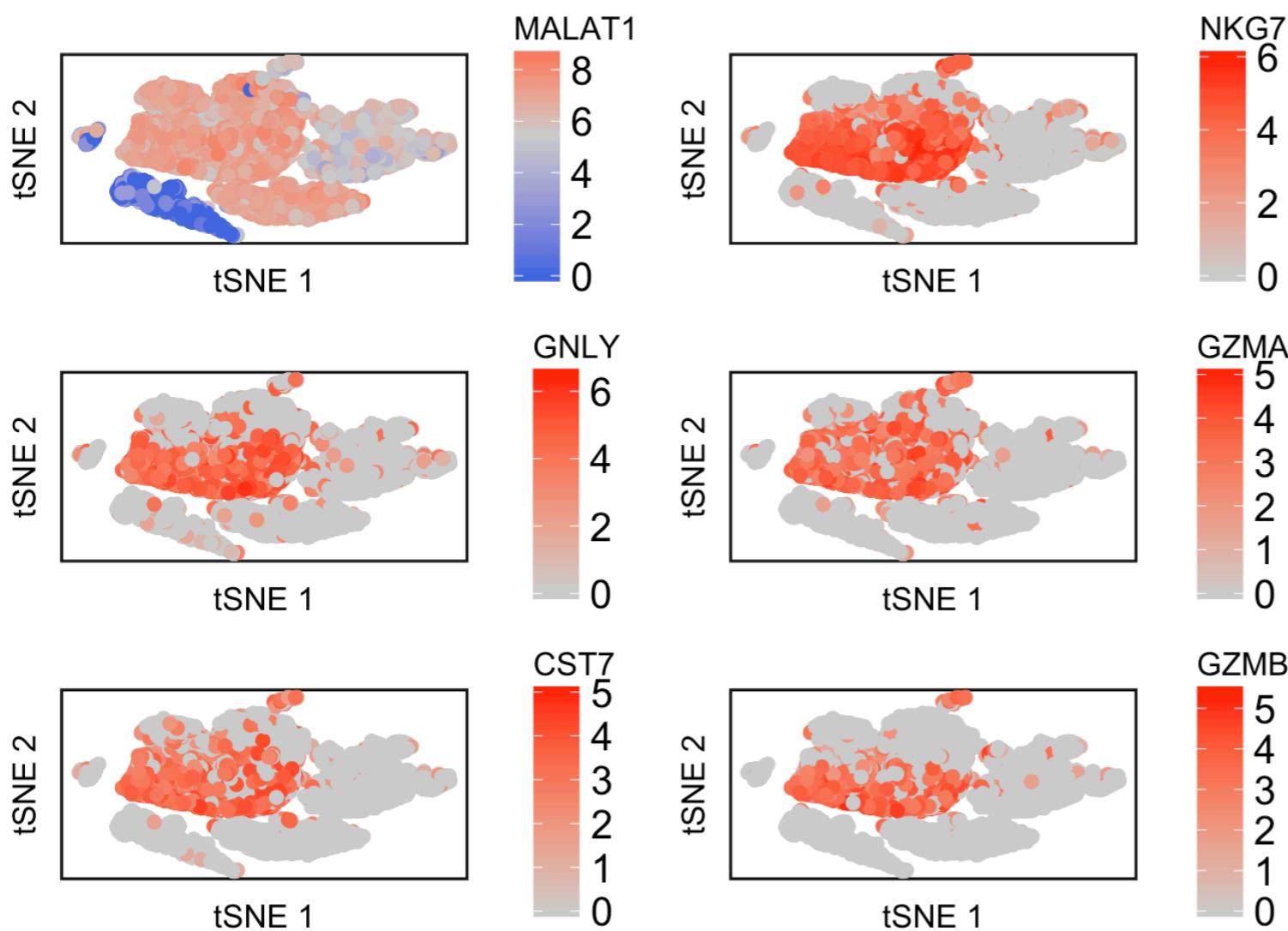
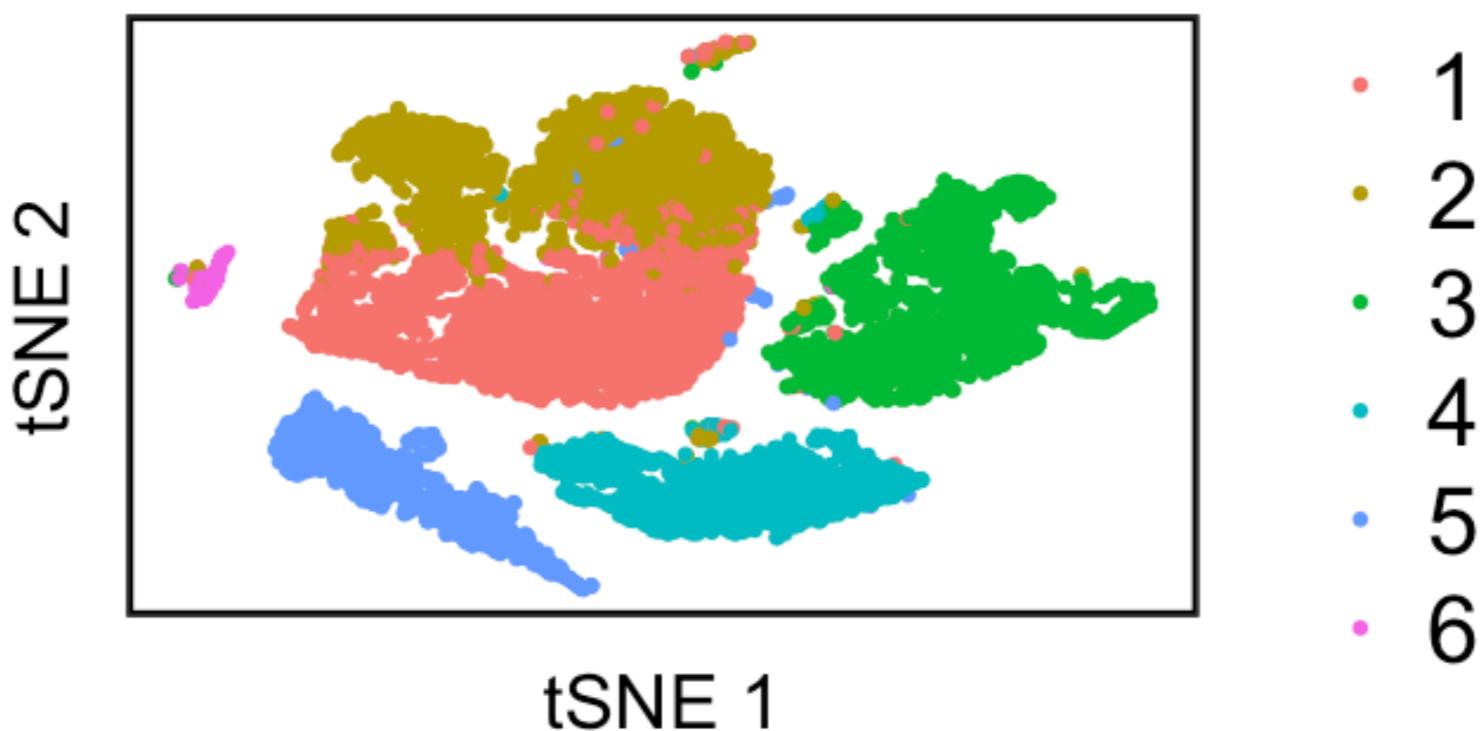
Time point:	Pre-ImmunoTx	Failing T cells & α -PD-1	Responding to T cells & α -PD-1 & α -CTLA	Failing T cells & α -PD-1 & α -CTLA
Samples for 10X-Genomics	Single-cell tumor digest (viable) TIL (viable) PBMC	PBMC	PBMC	Single-cell tumor digest (viable) TIL (viable) PBMC
Complementary Samples	BX: IHC. TIL, PBMC: T cell phenotype, Adaptive	PBMC: T cell phenotype, adaptive	BX: IHC TIL: Adaptive PBMC: T cell phenotype, adaptive	BX: IHC. TIL, PBMC: T cell phenotype, Adaptive

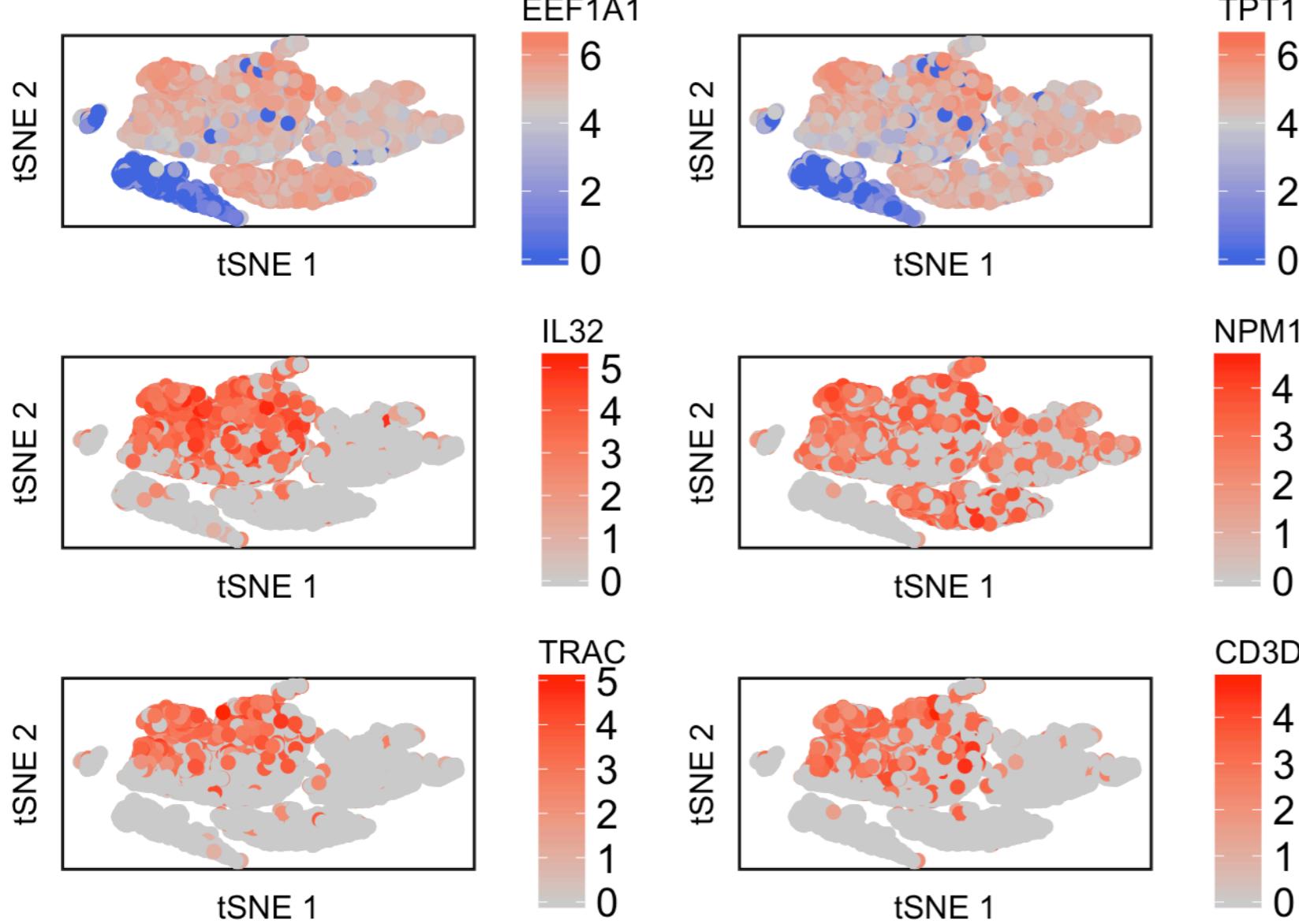
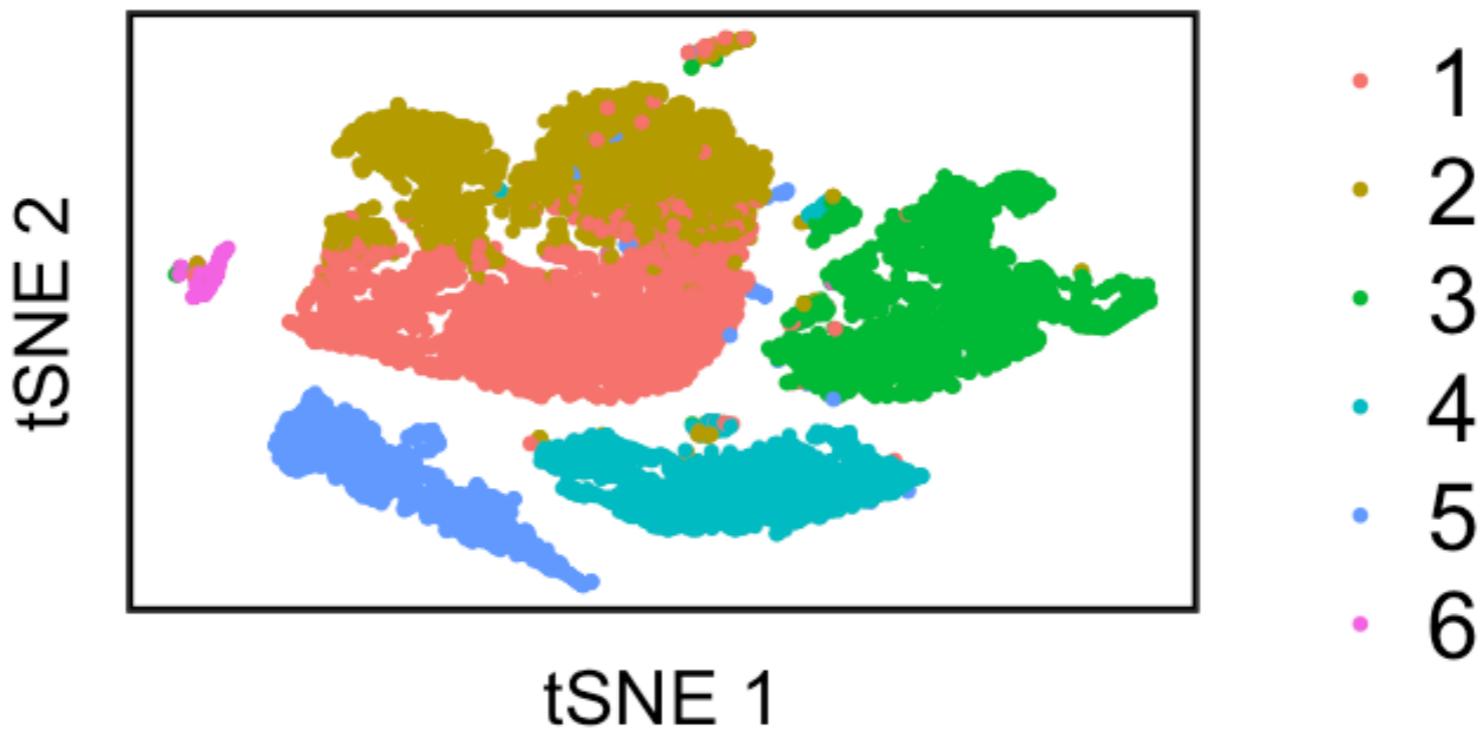




- Aug-2013
- Feb-2015
- Feb-2016
- Oct-2016

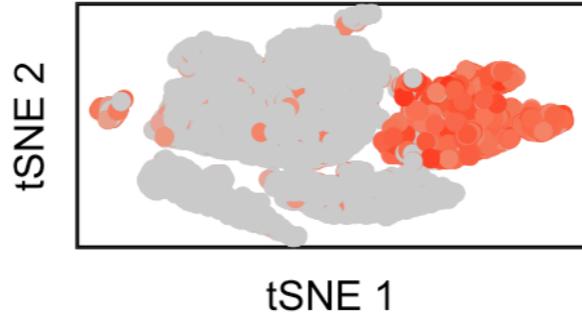
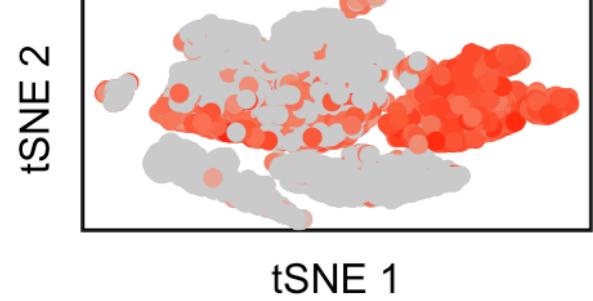
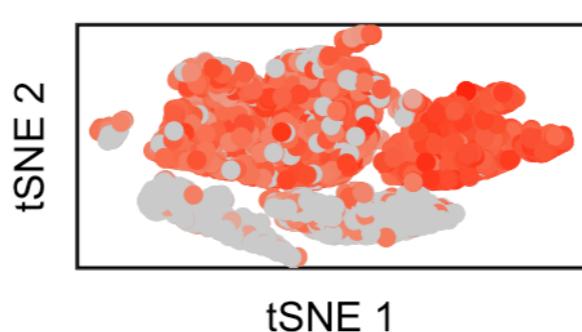
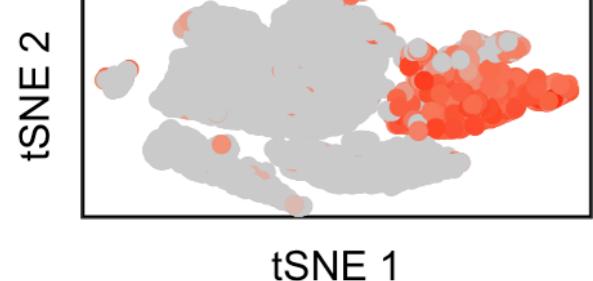
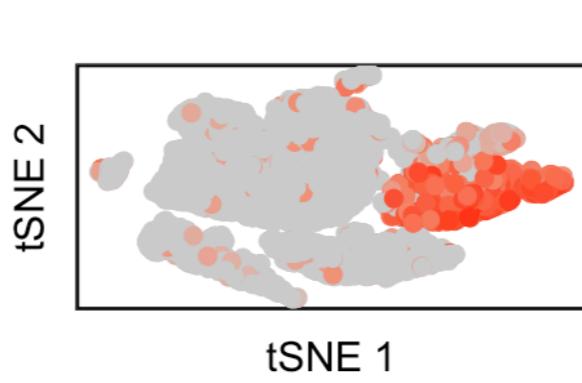
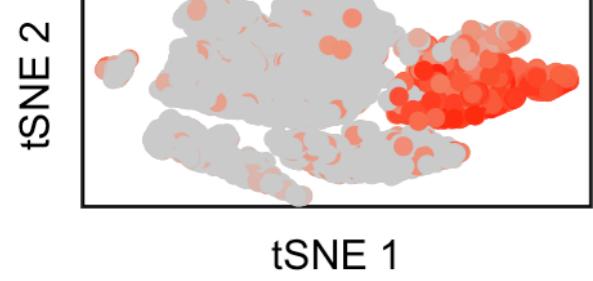
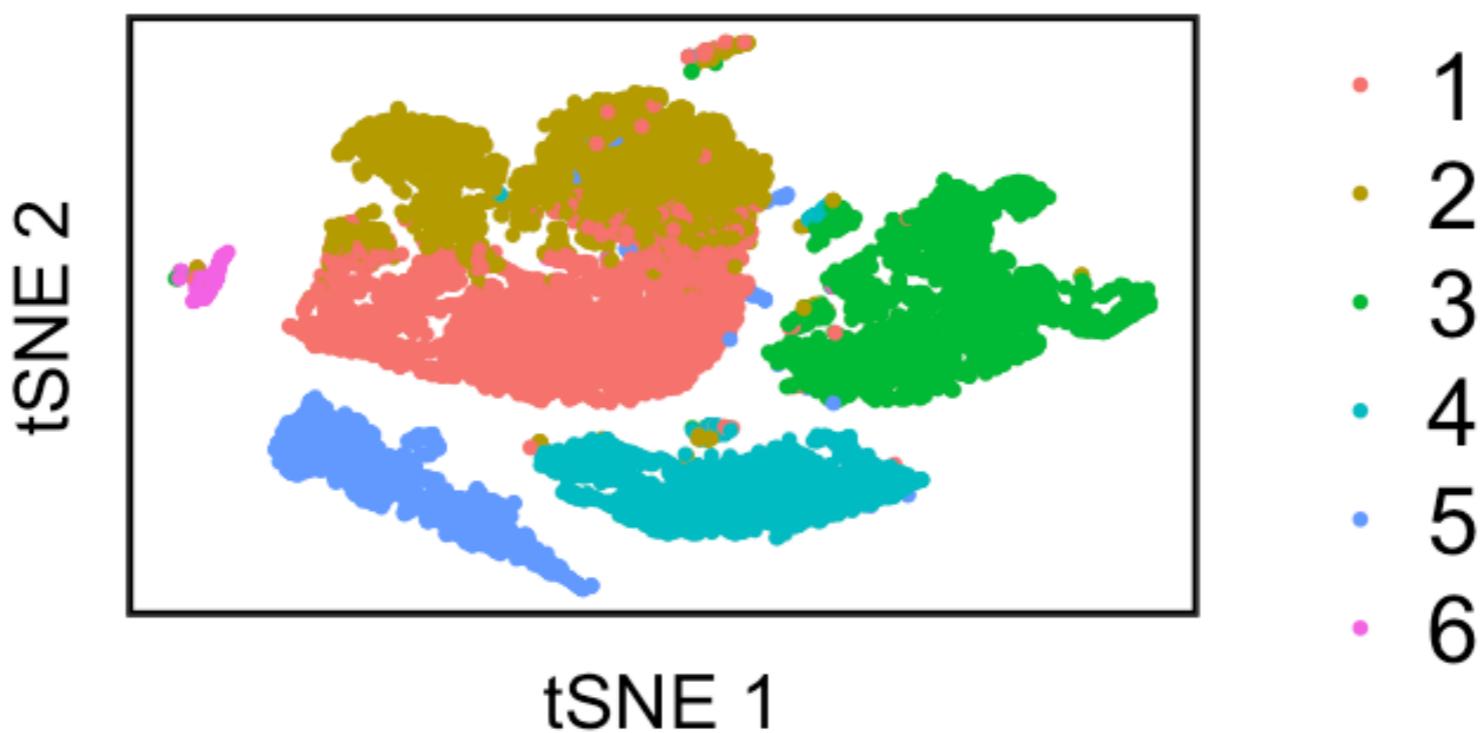






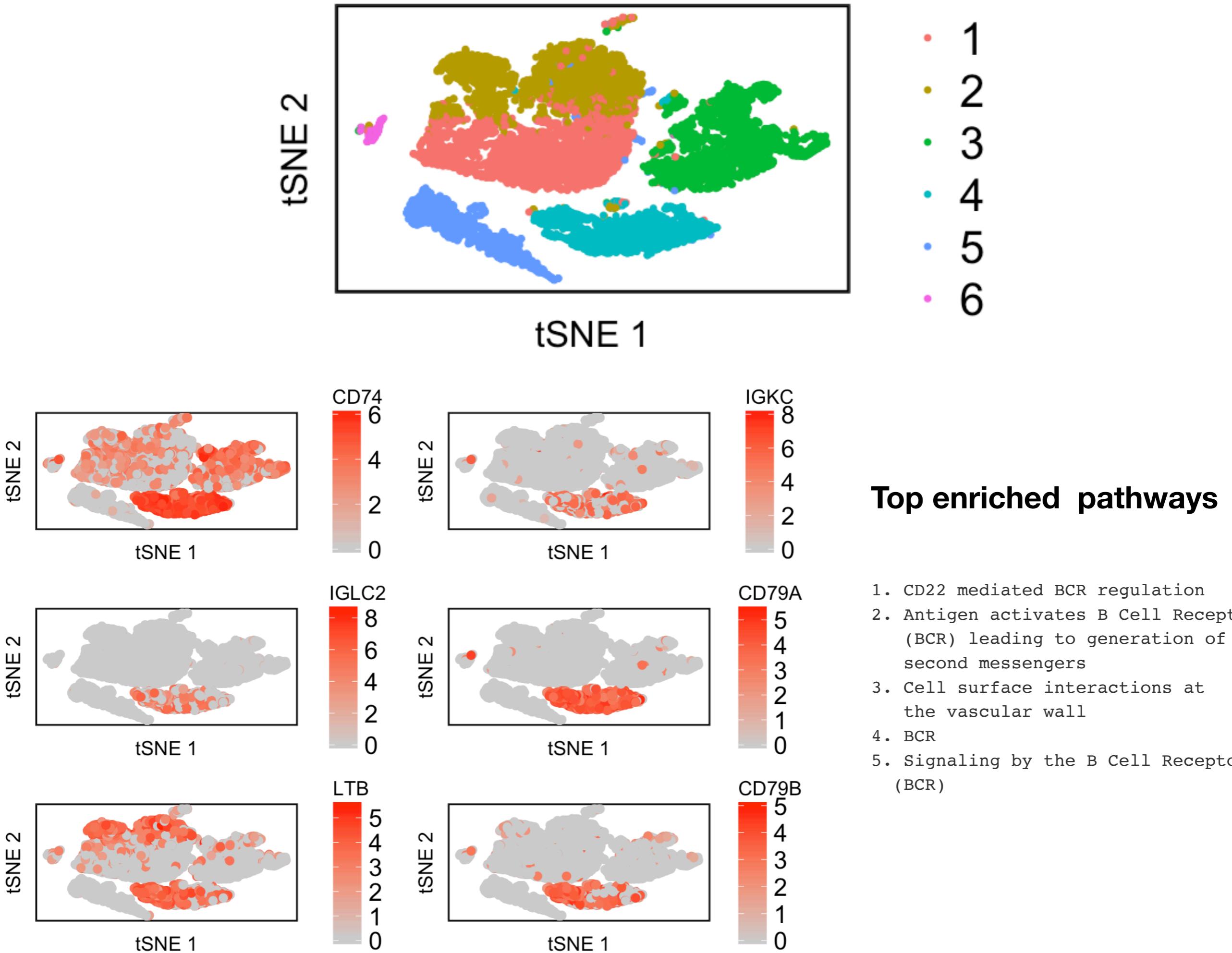
Top enriched pathways

1. Generation of second messenger molecules
2. Translocation of ZAP-70 to Immunological synapse
3. TCR
4. Phosphorylation of CD3 and TCR zeta chains
5. PD-1 signaling

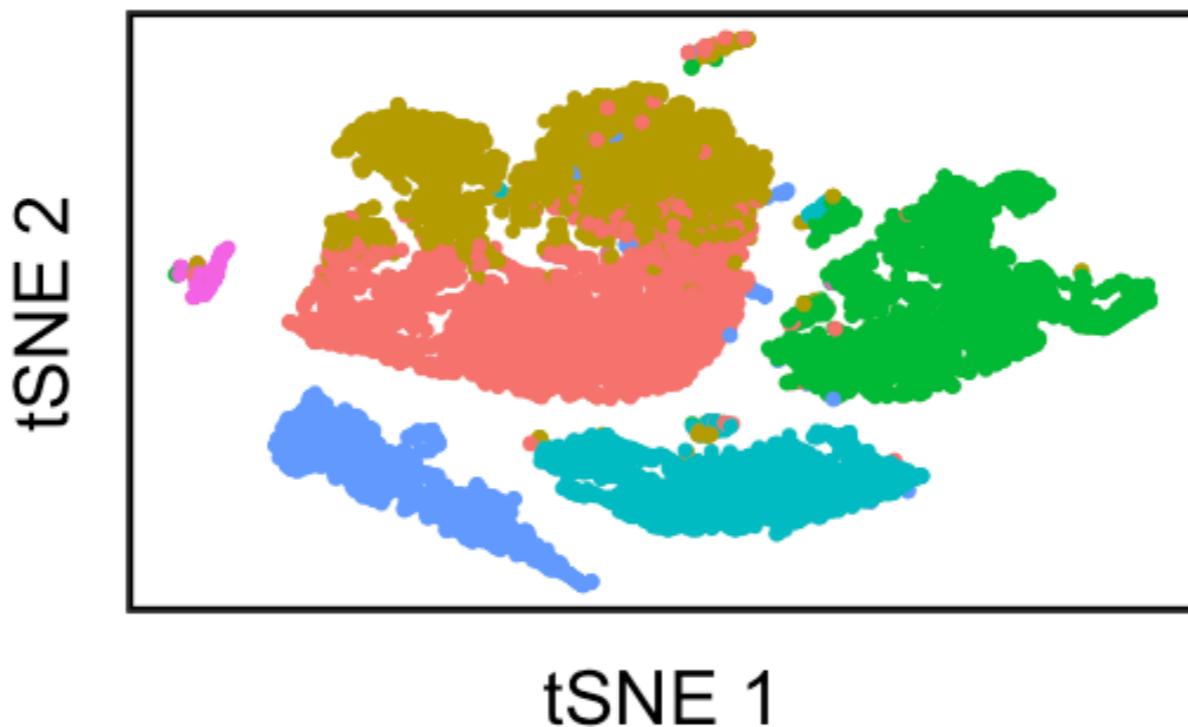


Top enriched pathways

1. Neutrophil degranulation
2. Innate Immune System
3. Immune System
4. Toll-Like Receptors Cascades
5. Endogenous TLR signaling
6. Regulation of TLR by endogenous ligand

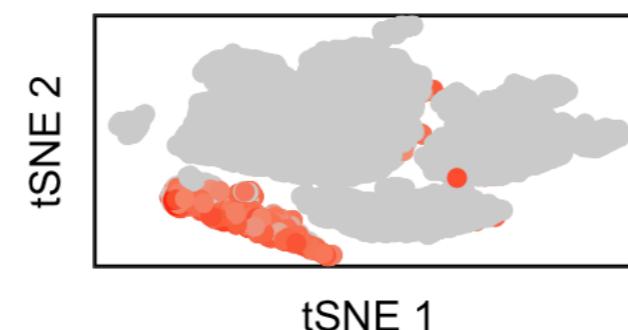
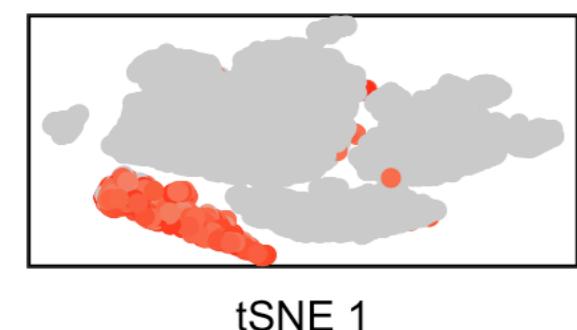
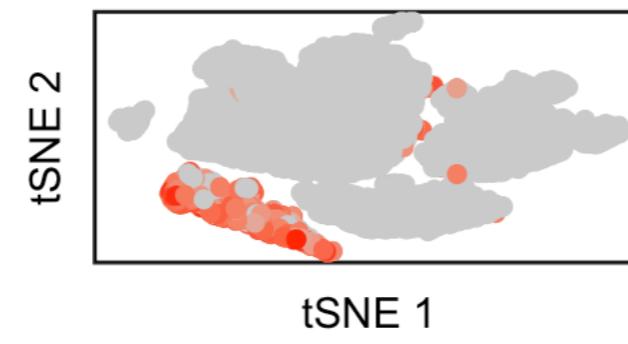
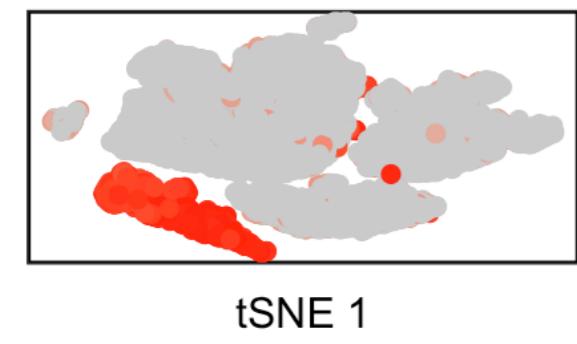
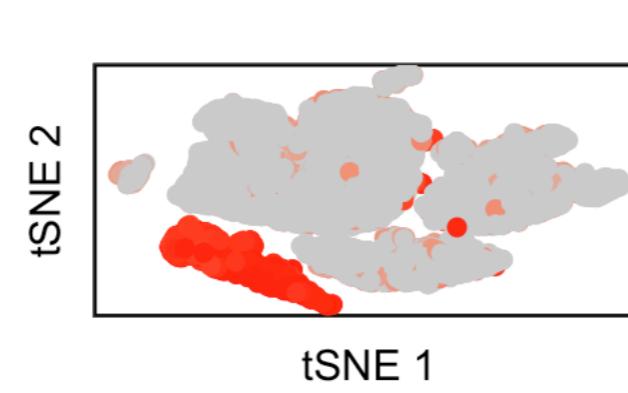
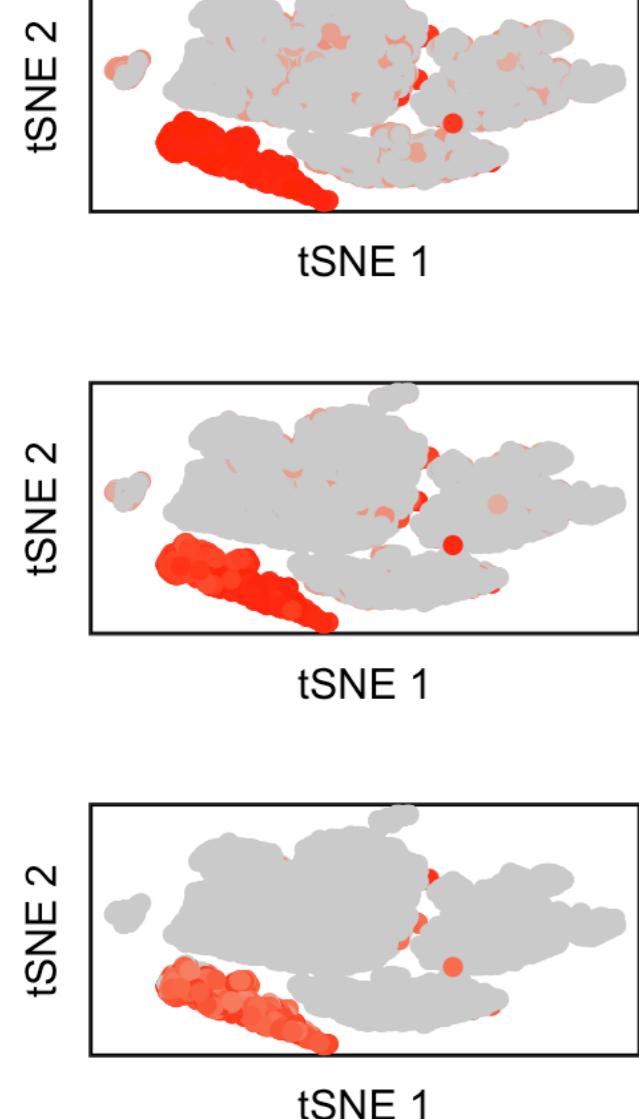


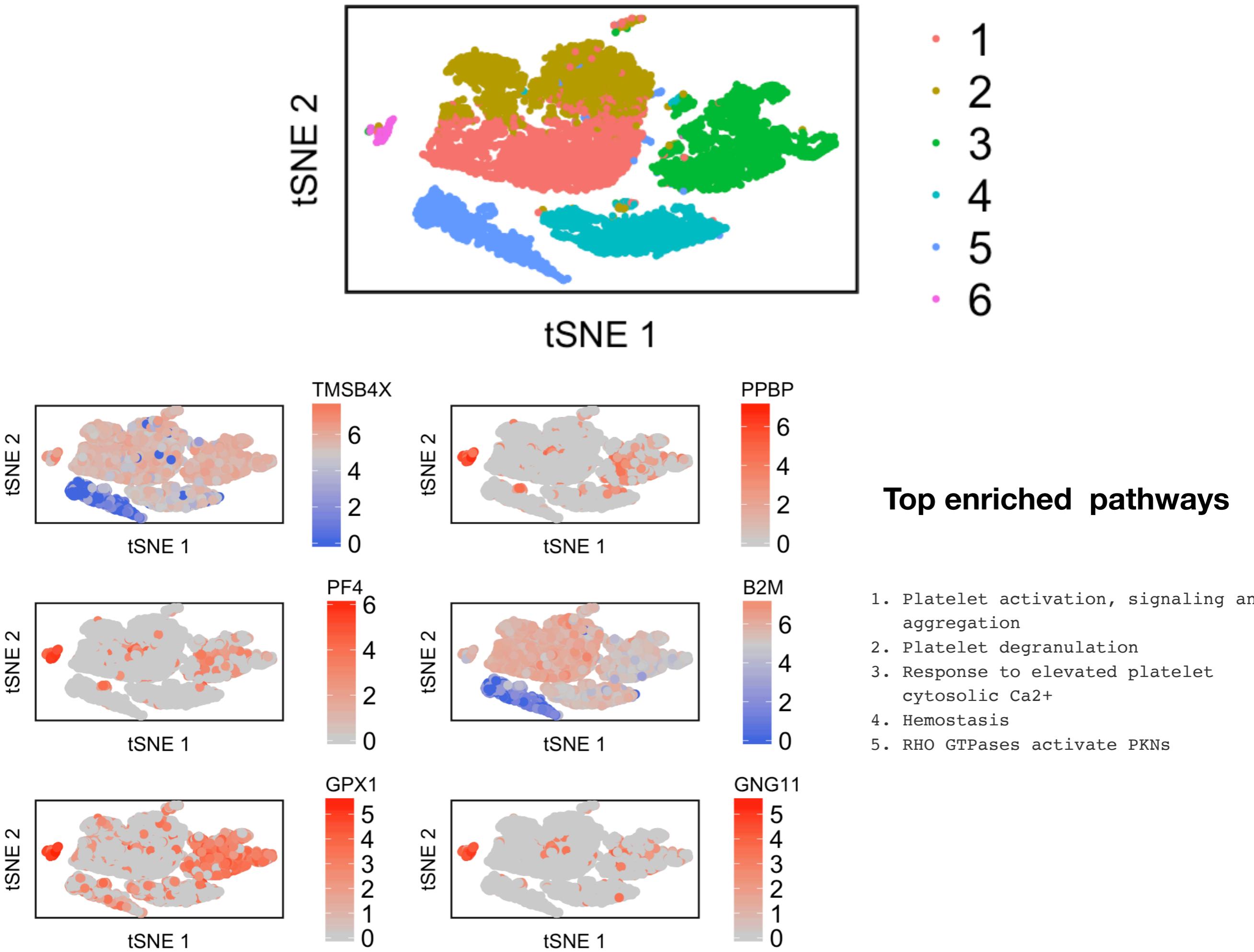
• 1
• 2
• 3
• 4
• 5
• 6



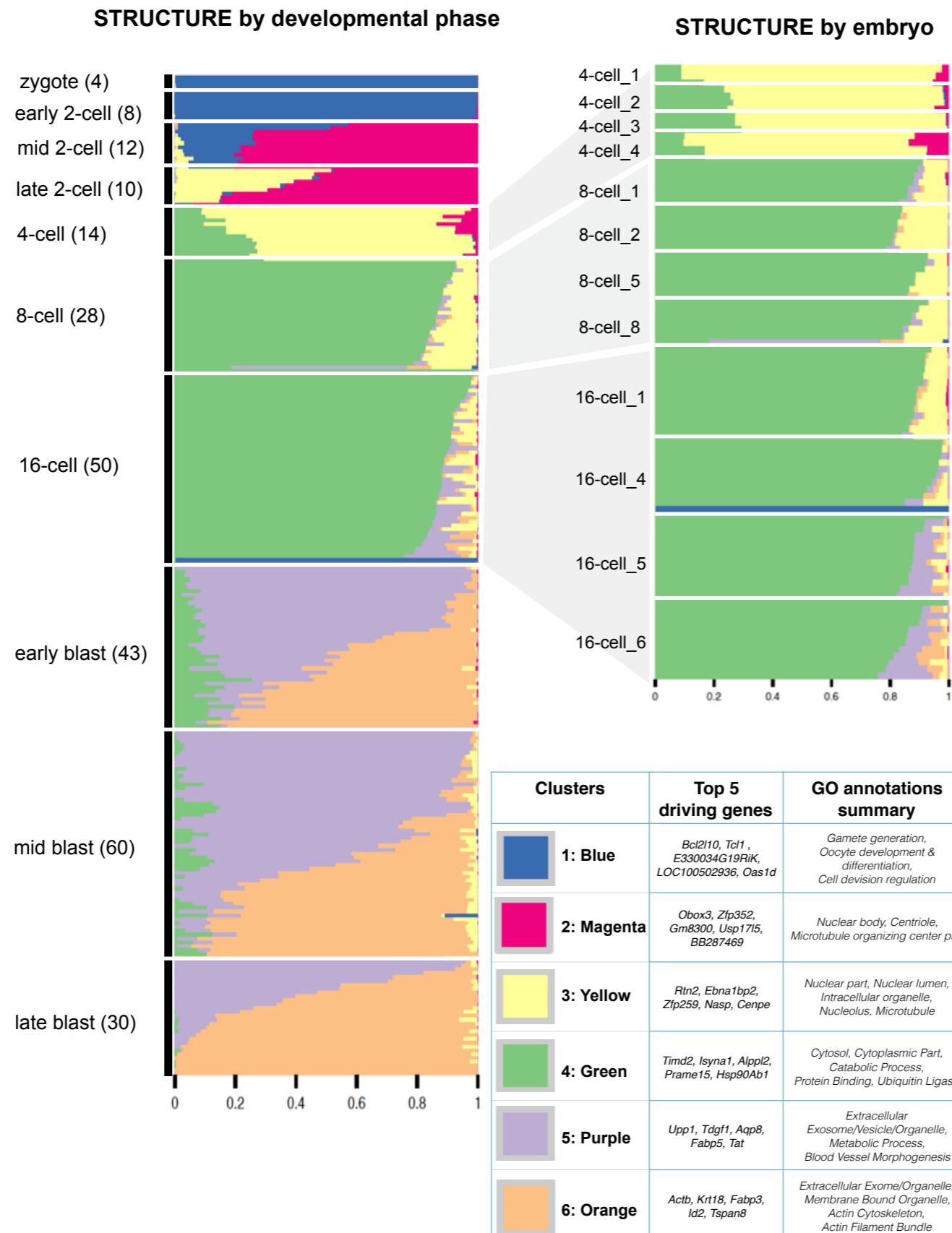
Top enriched pathways

1. hemoglobins chaperone
2. Erythrocytes take up oxygen and release carbon dioxide
3. Erythrocytes take up carbon dioxide and release oxygen
4. O₂/CO₂ exchange in erythrocytes
5. Malaria - Homo sapiens (human)





Deng et al 2016 single cell developmental study



Acknowledgements

**Matthew Stephens
Chiaowen Joyce Hsiao
Matt Taddy

GTEx Consortium**

**Raphael Gottardo
Valentin Voillet
Aude Chapuis
Kelly Paulson
Paul Ngheim**

Project webpage:

<https://kkdey.github.io/singlecell-clustering/>

<http://stephenslab.github.io/count-clustering/>

Package : CountClust

devel: <https://github.com/kkdey/CountClust>