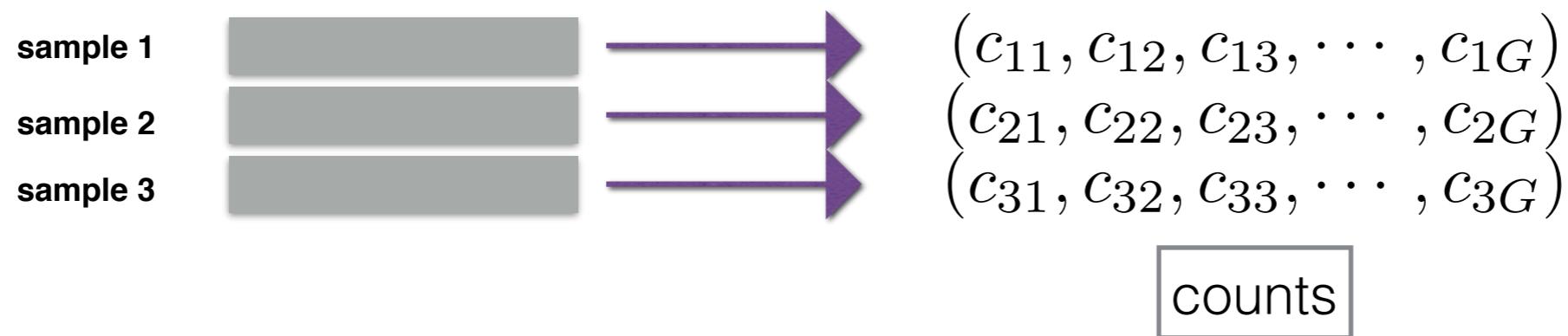


Model-Based Clustering and Visualization of RNA-Seq Data

Kushal Dey
Joyce Hsiao
Matthew Stephens

RNA-seq read counts data



- We consider number of read counts per gene
- Samples may come from bulk RNA-seq (tissue samples) or single cell RNA-seq (scRNA-seq).

Bulk RNA-seq



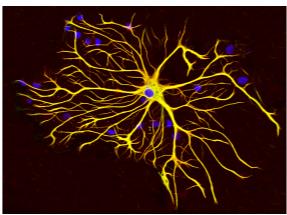
cell type 1

example:
brain



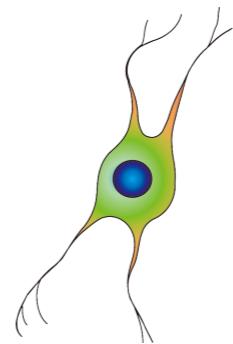
Neuron

cell type 2



Astrocyte

cell type 3



Oligodendrocytes



cell type 1

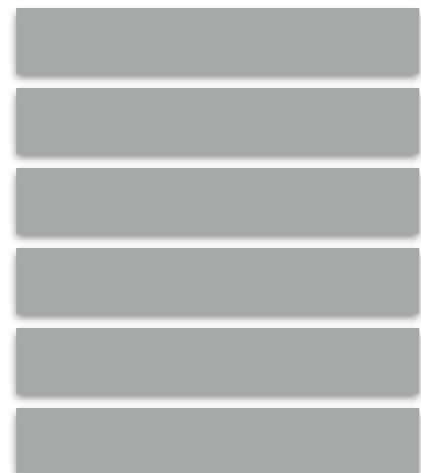
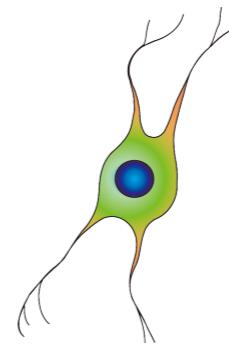
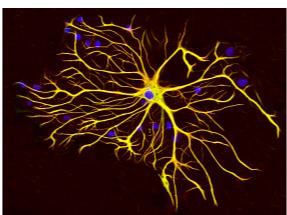
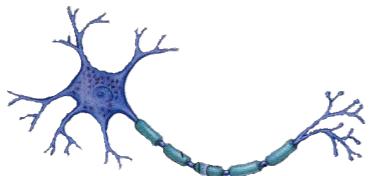


cell type 2

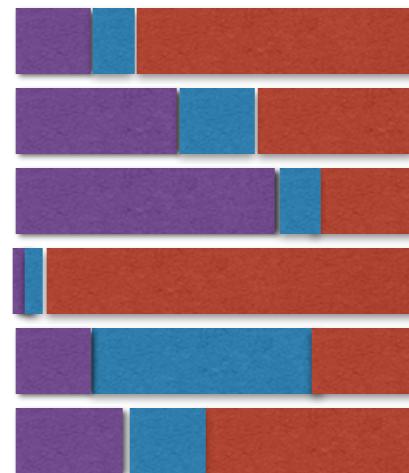
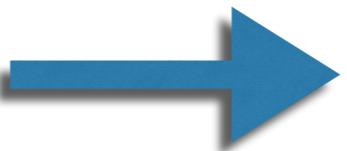


cell type 3

example:
brain

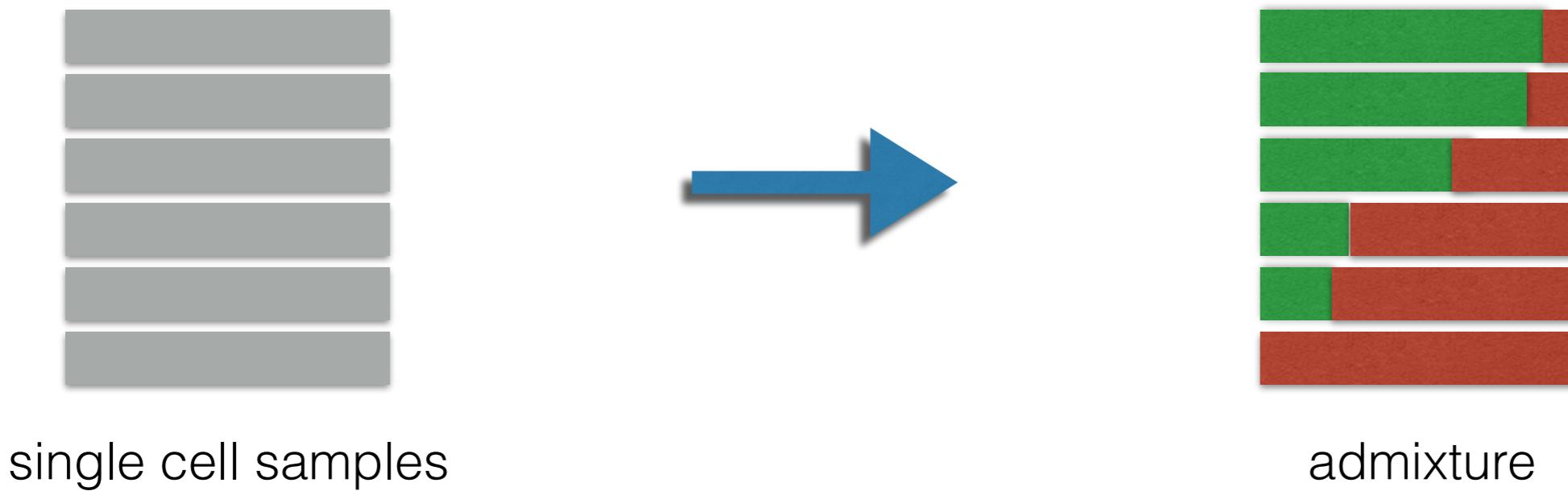
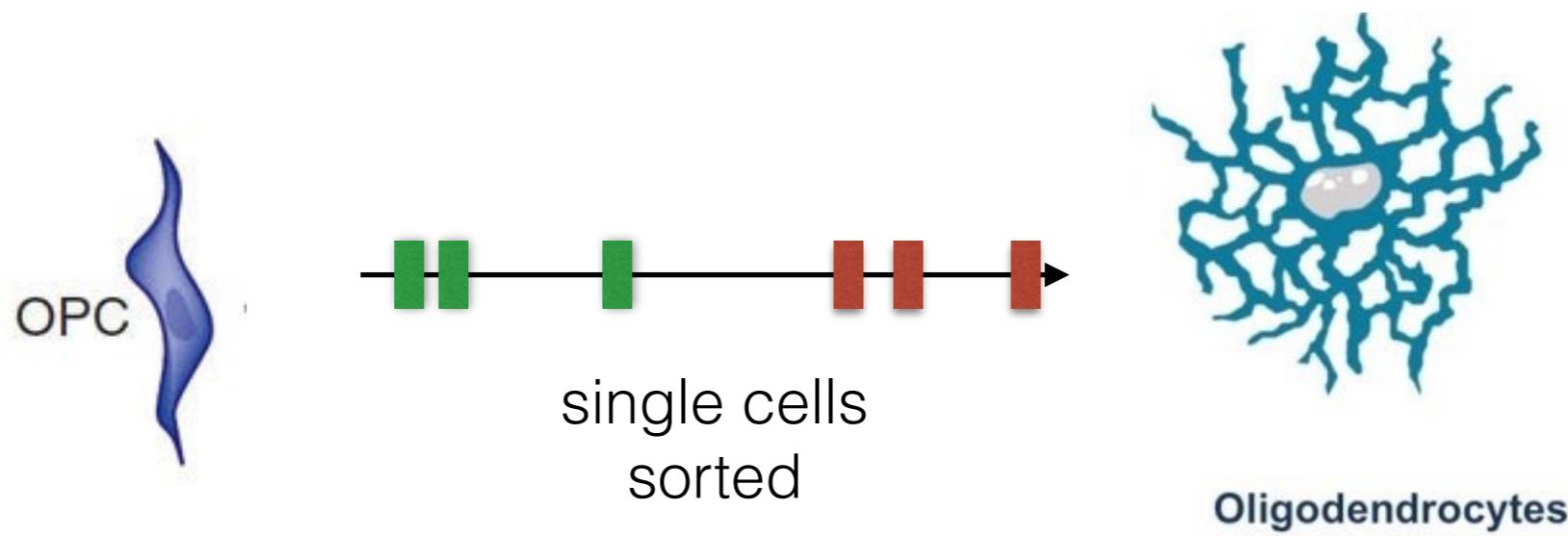


bulk RNA
tissue samples



admixture

Single cell RNA-seq

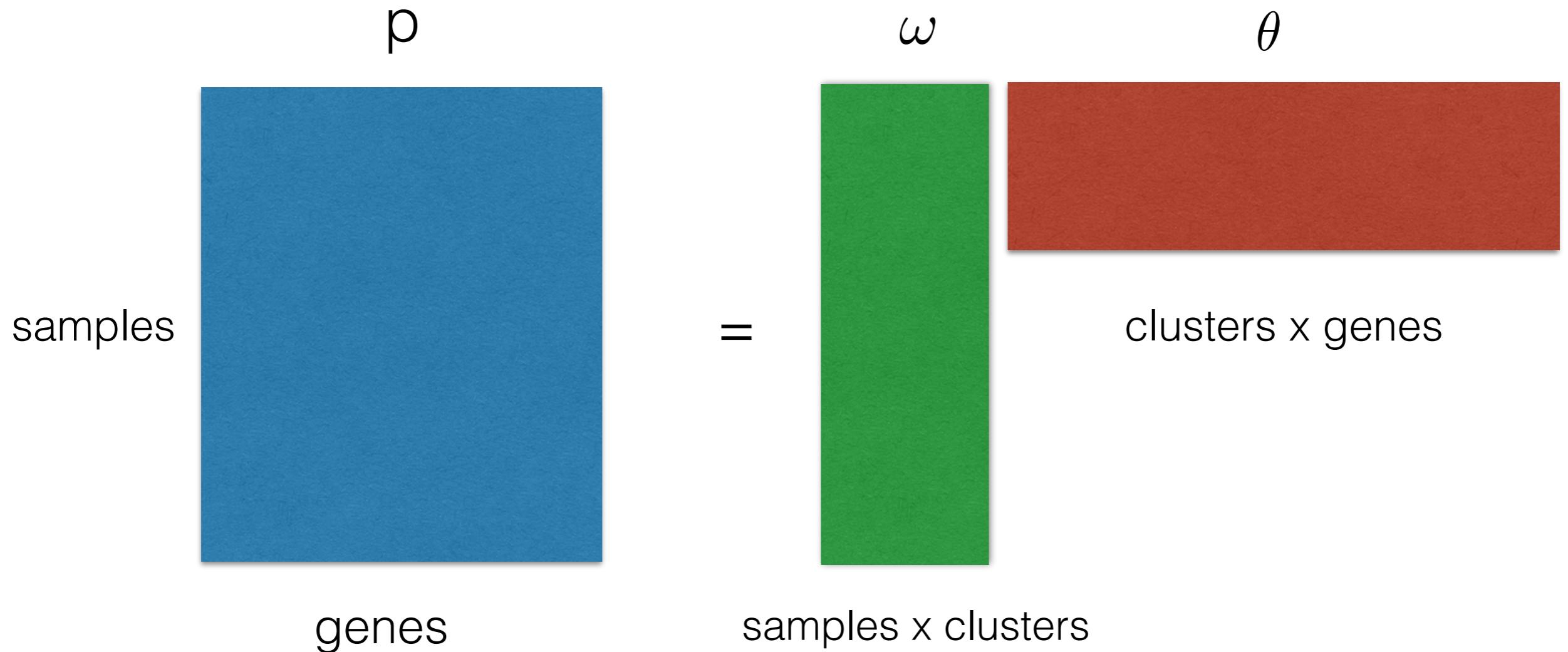


Model

Grade of Membership (GoM) Model

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

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



Grade of Membership Model (continued)

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

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

ω_{nk} grade of membership of sample n in cluster or topic k

θ_{kg} relative proportion of expression of gene g in cluster k

ω_{nk} and θ_{kg} were determined in an unsupervised manner
assuming Dirichlet priors

$$\omega_{n*} \sim Dir\left(\frac{1}{K}, \frac{1}{K}, \dots, \frac{1}{K}\right) \quad \theta_{k*} \sim Dir\left(\frac{1}{KG}, \frac{1}{KG}, \dots, \frac{1}{KG}\right)$$

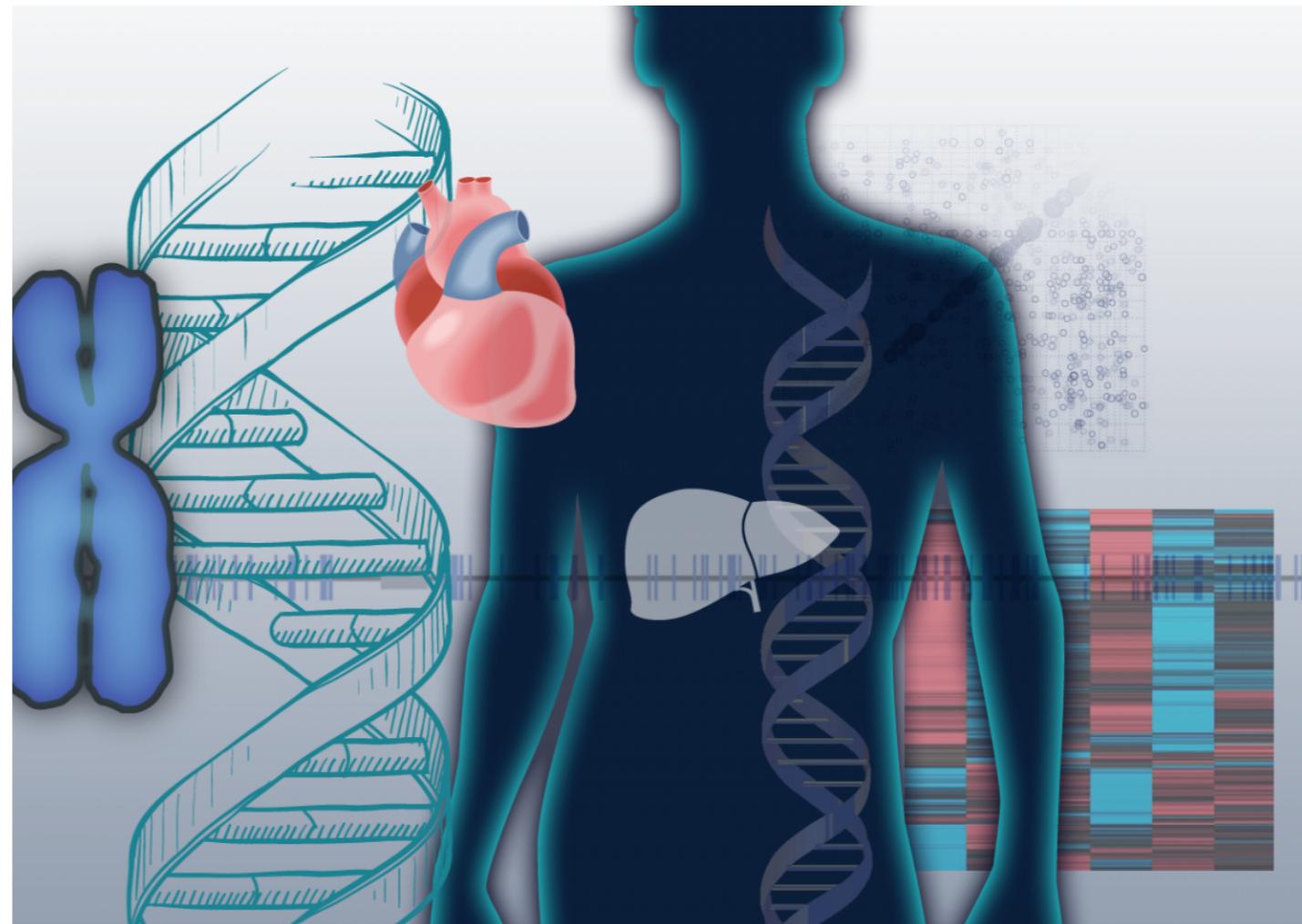
Data Analysis

Genotype Tissue Expression Project - V6

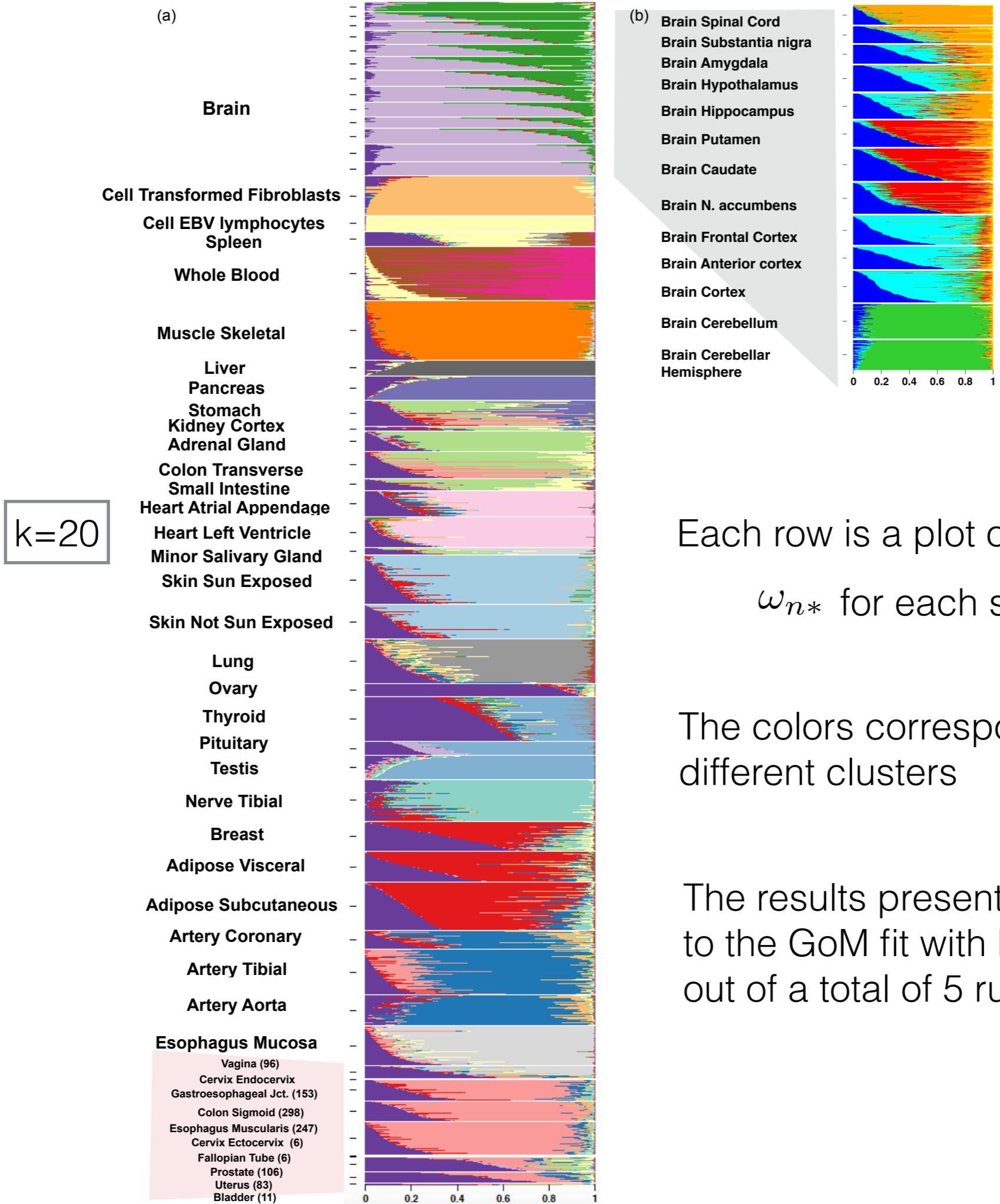
has data on 8555 tissue samples taken from 900 donors (Phase 6) for around 56,000 genes.

The tissue samples came from a total of 53 different tissues/ sub-tissues (e.g. - brain sub tissues, whole blood etc)

We filtered out around 16,069 genes that passed QC tests.



We fit the GoM model and get estimates of ω and θ



Each row is a plot of the vector
 ω_{n*} for each sample n

The colors correspond to
different clusters

The results presented is corresponding
to the GoM fit with highest Bayes factor
out of a total of 5 runs.

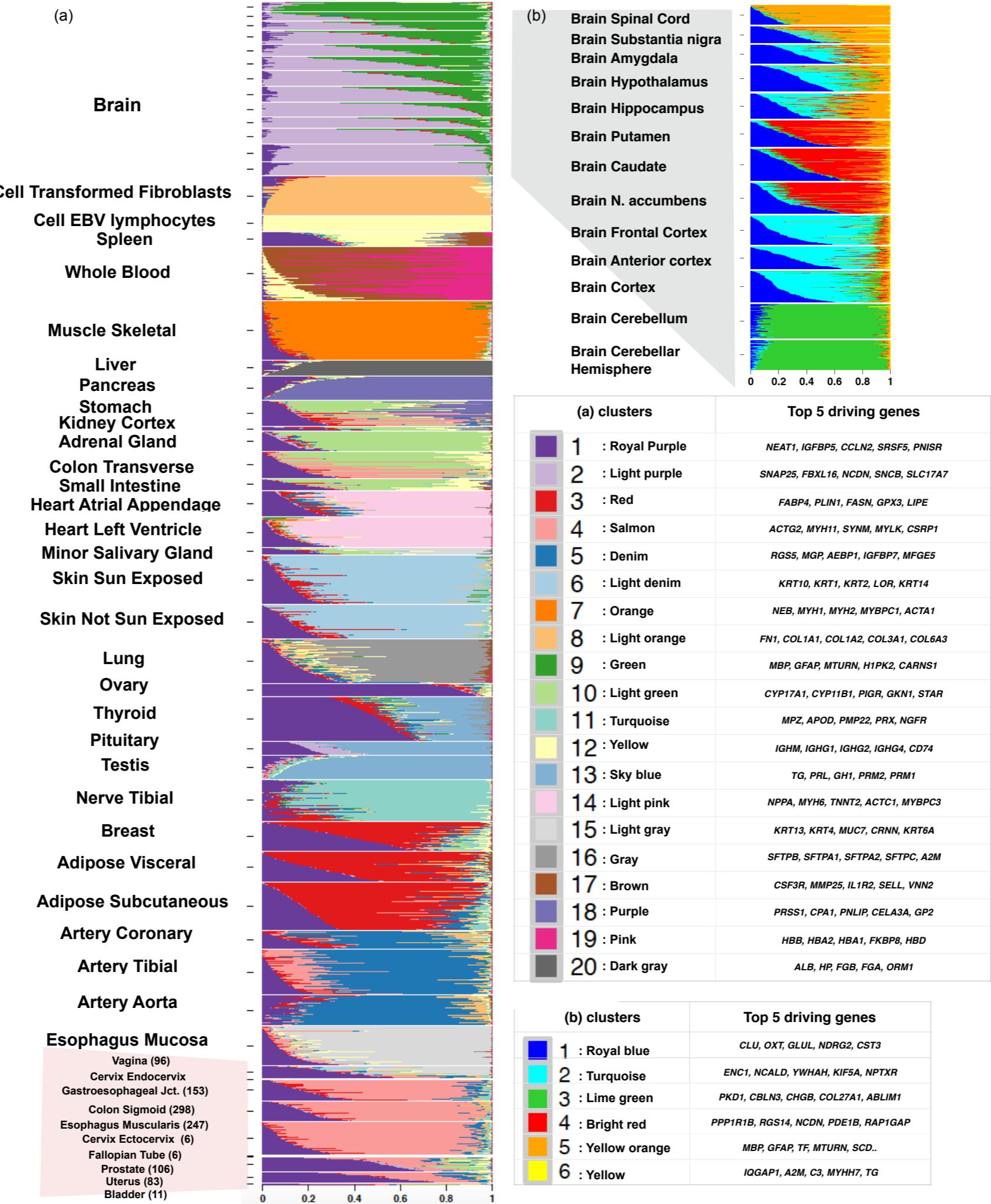
Extracting top genes per cluster

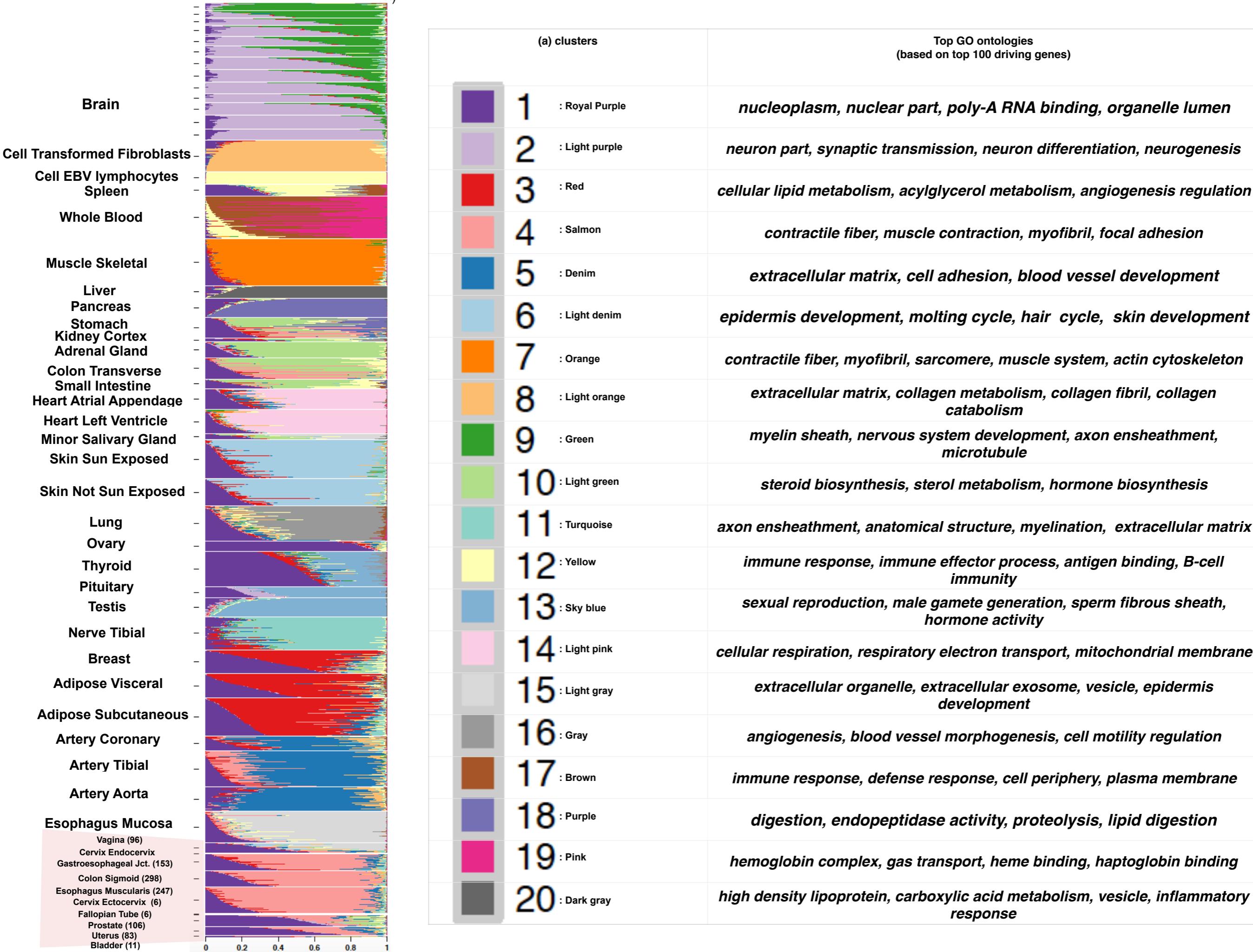
We compute for each cluster k and each gene g

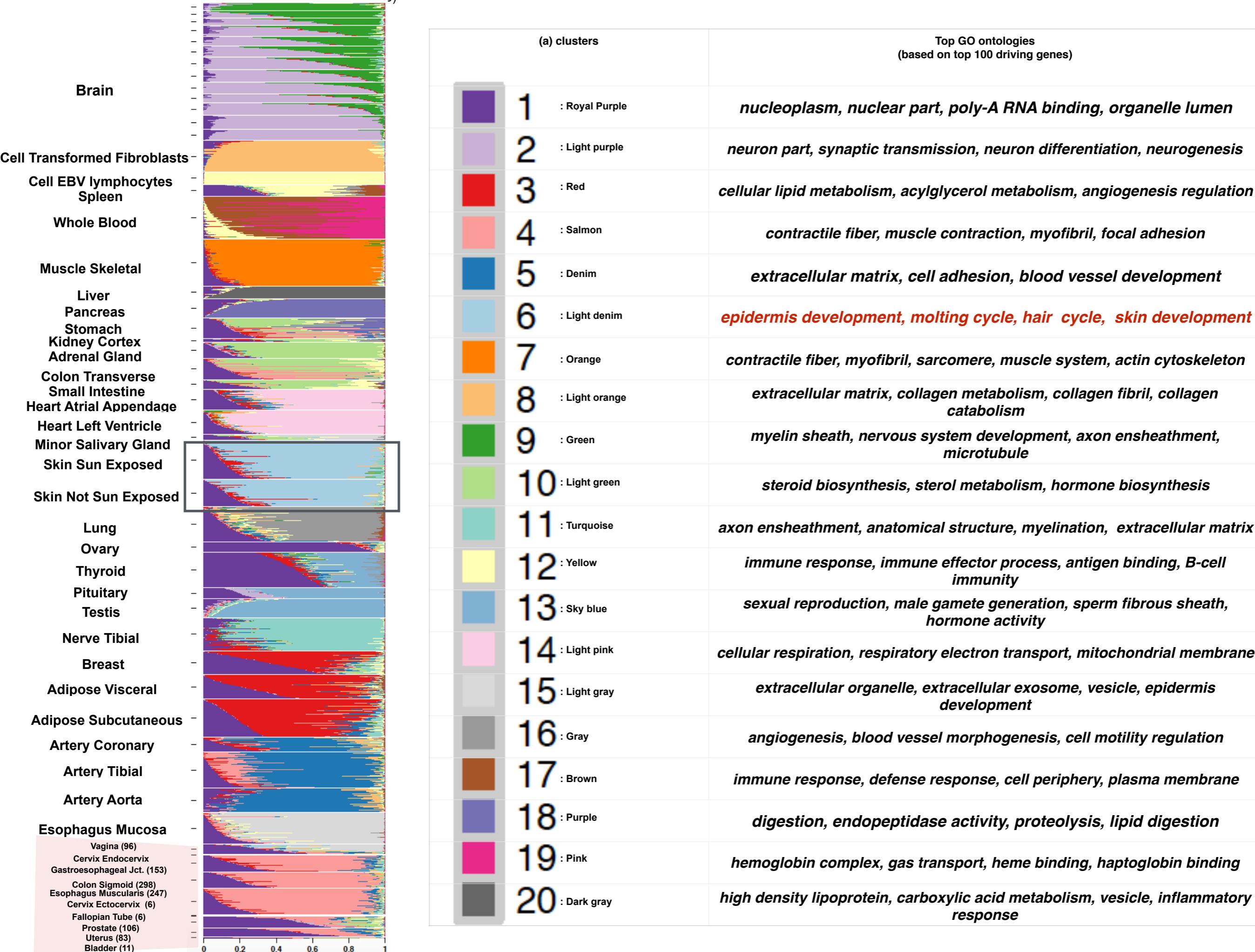
$$D^g[k] = \min_{l \neq k} \left[\hat{\theta}_{kg} \log \left(\frac{\hat{\theta}_{kg}}{\hat{\theta}_{lg}} \right) + \hat{\theta}_{lg} - \hat{\theta}_{kg} \right]$$

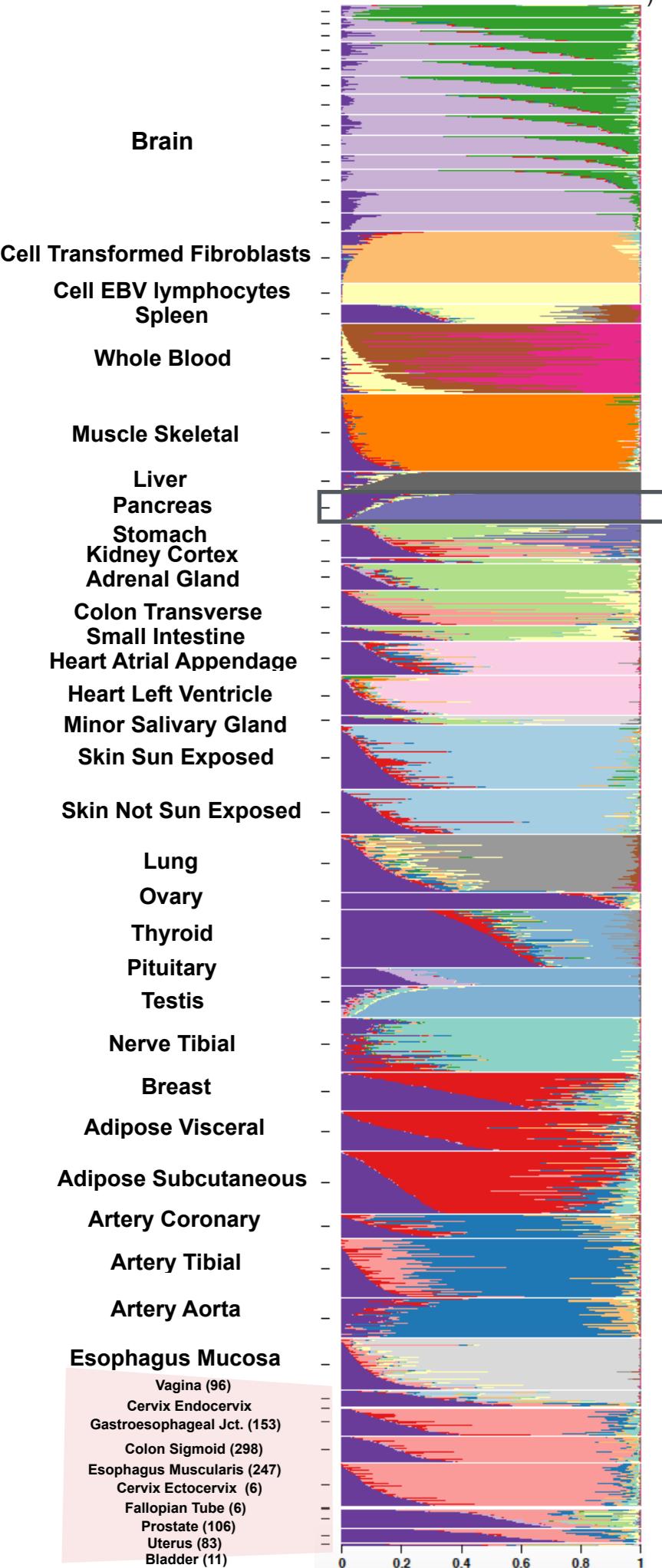
For each k , choose the genes with the highest values of $D^g[k]$

which satisfy $\hat{\theta}_{kg} = \max_l \hat{\theta}_{lg}$

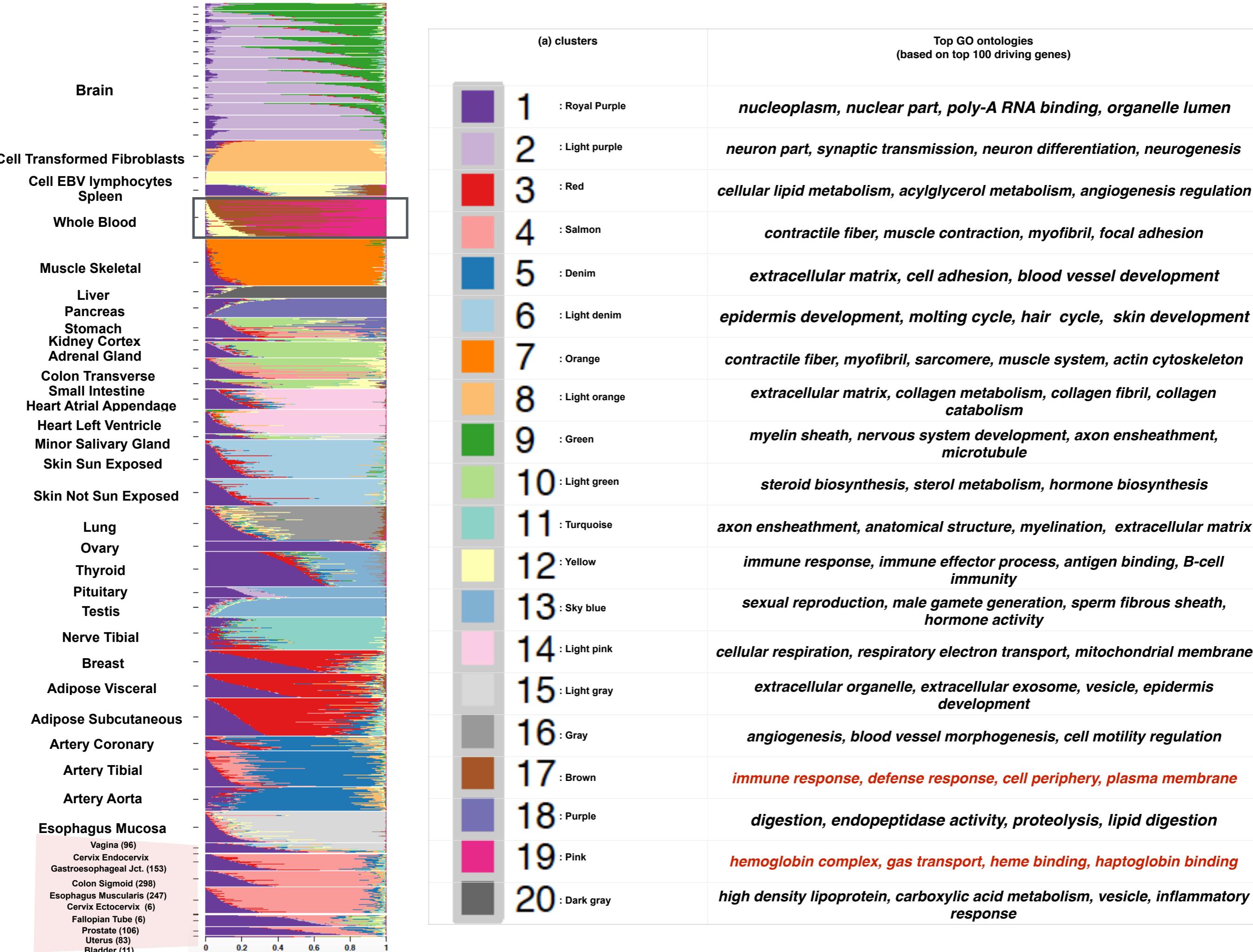


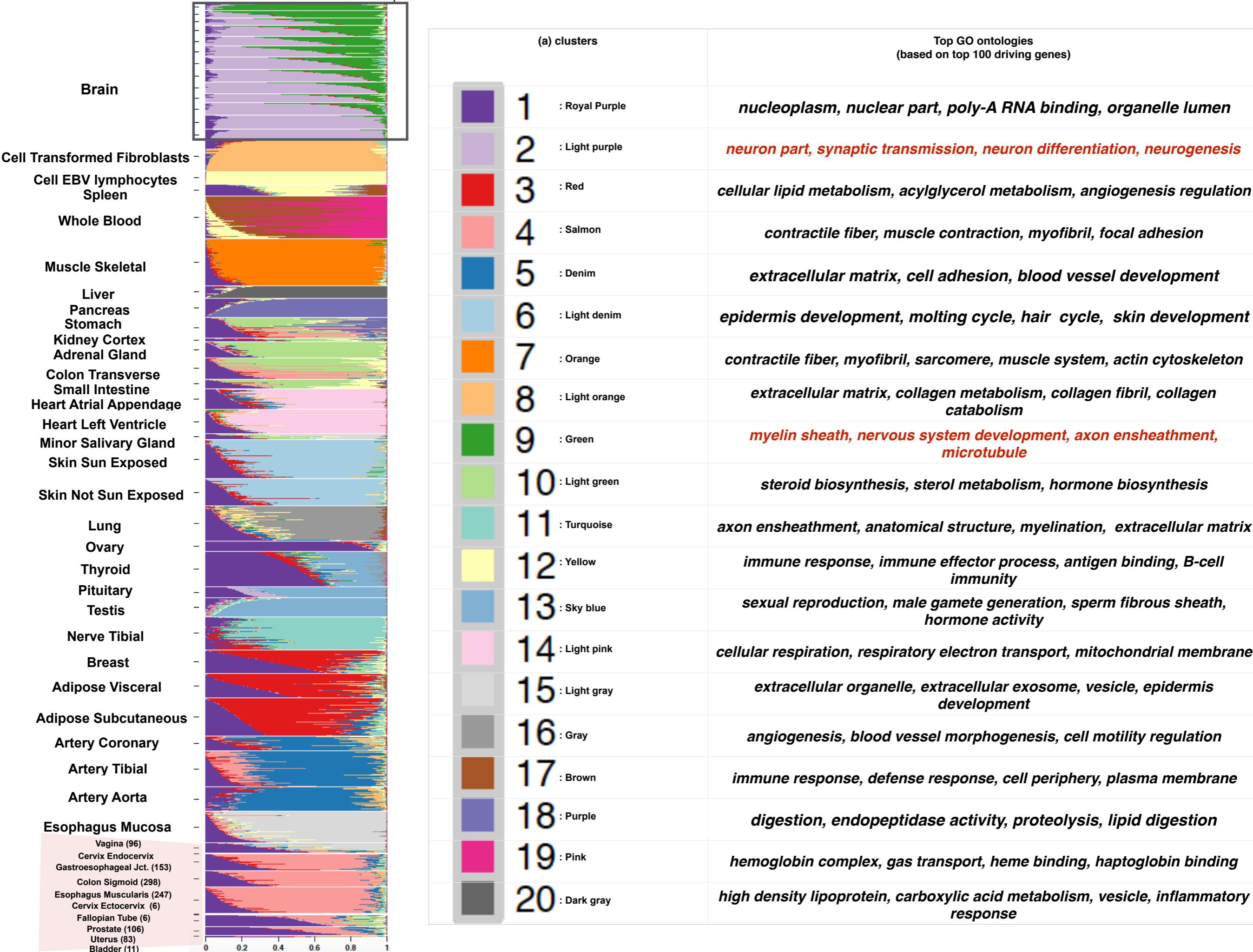


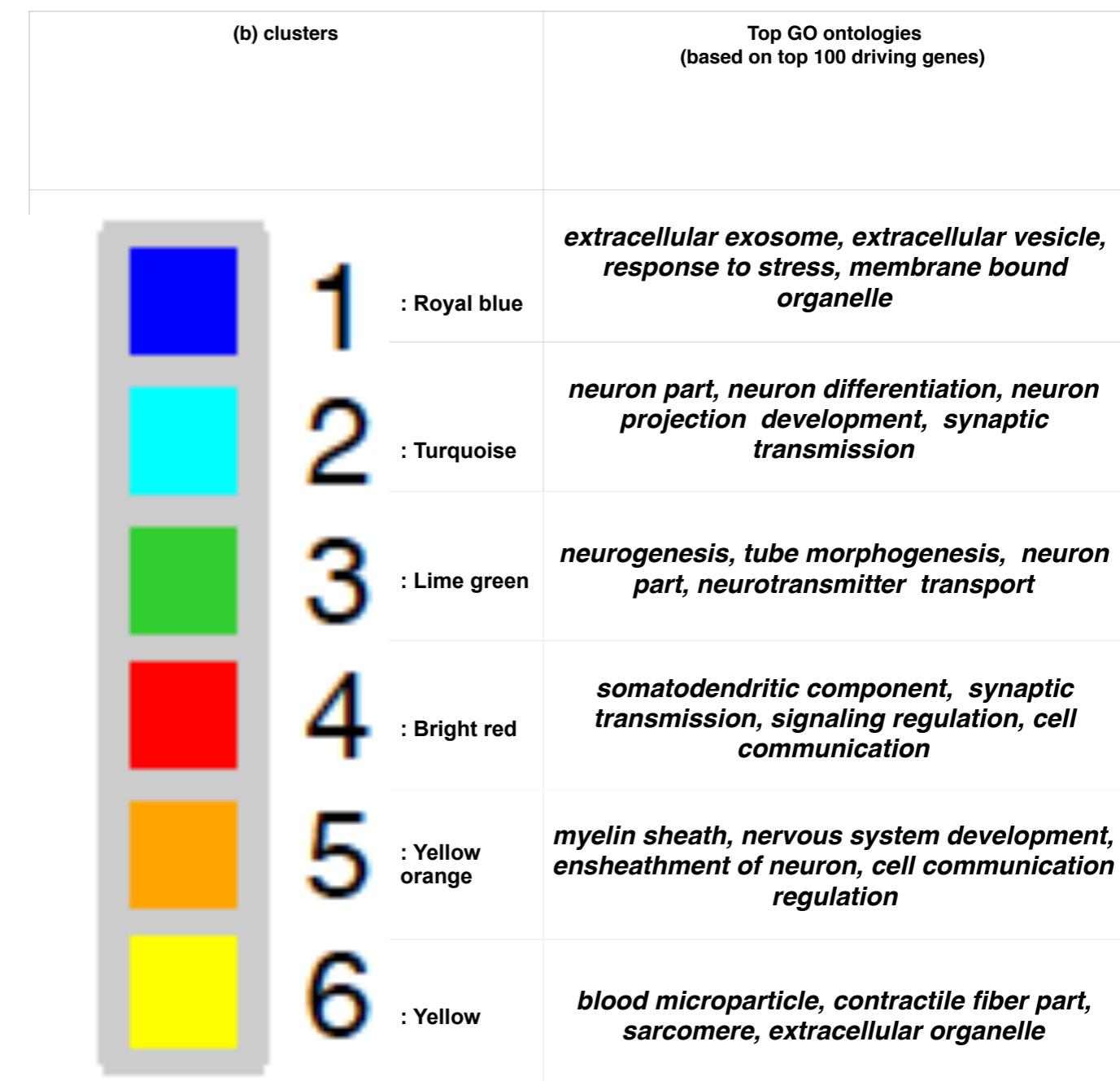
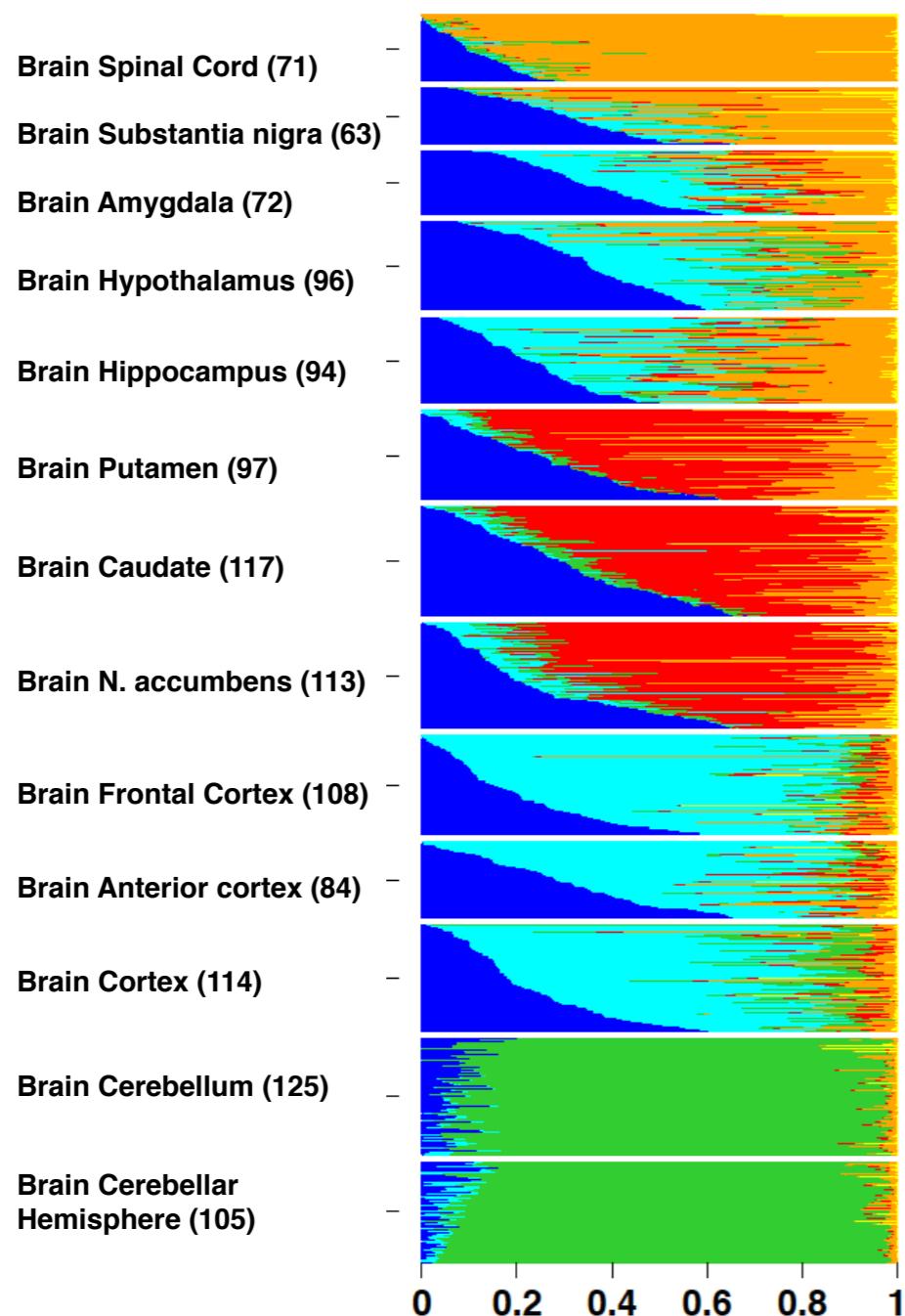


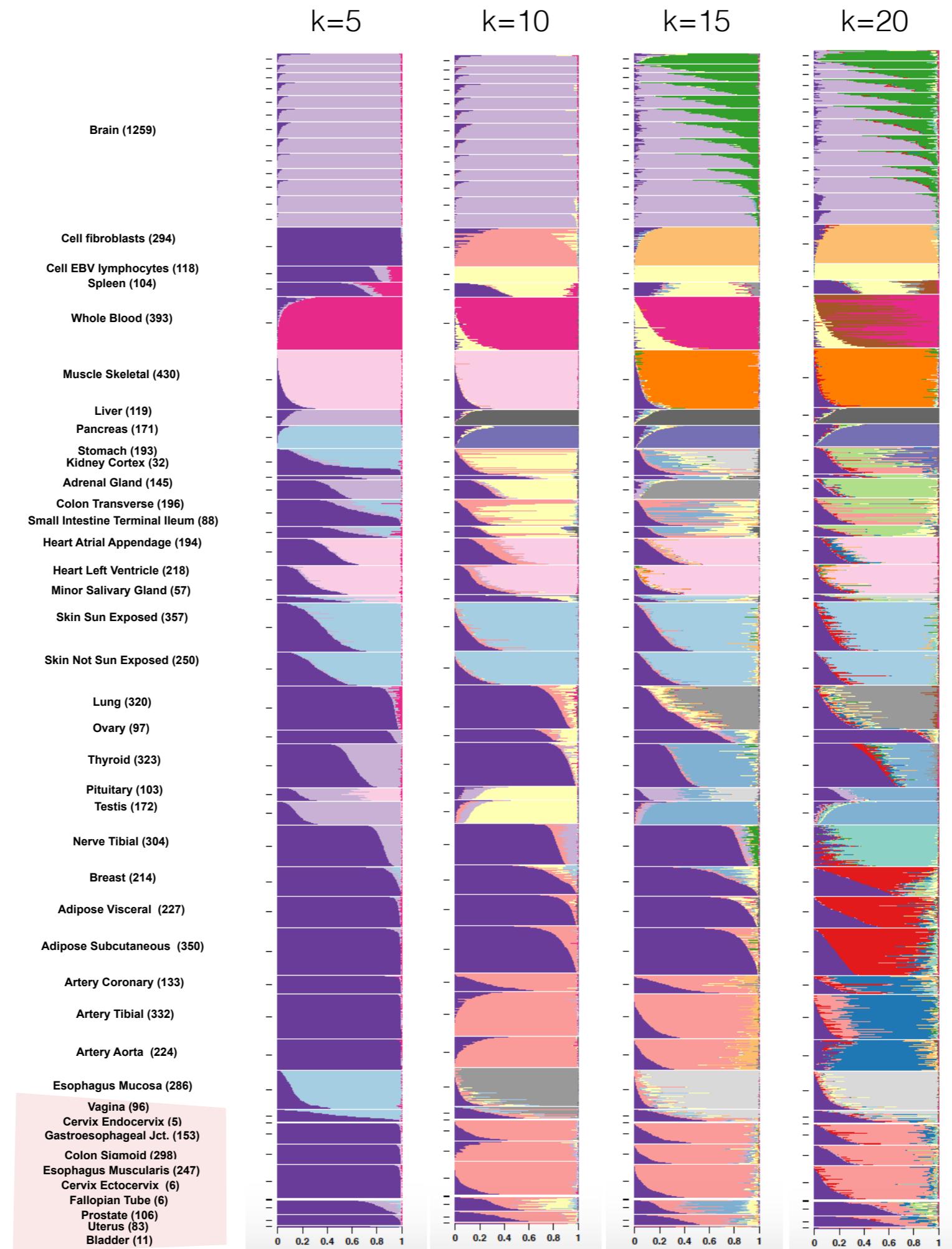


(a) clusters	Top GO ontologies (based on top 100 driving genes)
1 : Royal Purple	nucleoplasm, nuclear part, poly-A RNA binding, organelle lumen
2 : Light purple	neuron part, synaptic transmission, neuron differentiation, neurogenesis
3 : Red	cellular lipid metabolism, acylglycerol metabolism, angiogenesis regulation
4 : Salmon	contractile fiber, muscle contraction, myofibril, focal adhesion
5 : Denim	extracellular matrix, cell adhesion, blood vessel development
6 : Light denim	epidermis development, molting cycle, hair cycle, skin development
7 : Orange	contractile fiber, myofibril, sarcomere, muscle system, actin cytoskeleton
8 : Light orange	extracellular matrix, collagen metabolism, collagen fibril, collagen catabolism
9 : Green	myelin sheath, nervous system development, axon ensheathment, microtubule
10 : Light green	steroid biosynthesis, sterol metabolism, hormone biosynthesis
11 : Turquoise	axon ensheathment, anatomical structure, myelination, extracellular matrix
12 : Yellow	immune response, immune effector process, antigen binding, B-cell immunity
13 : Sky blue	sexual reproduction, male gamete generation, sperm fibrous sheath, hormone activity
14 : Light pink	cellular respiration, respiratory electron transport, mitochondrial membrane
15 : Light gray	extracellular organelle, extracellular exosome, vesicle, epidermis development
16 : Gray	angiogenesis, blood vessel morphogenesis, cell motility regulation
17 : Brown	immune response, defense response, cell periphery, plasma membrane
18 : Purple	<i>digestion, endopeptidase activity, proteolysis, lipid digestion</i>
19 : Pink	<i>hemoglobin complex, gas transport, heme binding, haptoglobin binding</i>
20 : Dark gray	<i>high density lipoprotein, carboxylic acid metabolism, vesicle, inflammatory response</i>







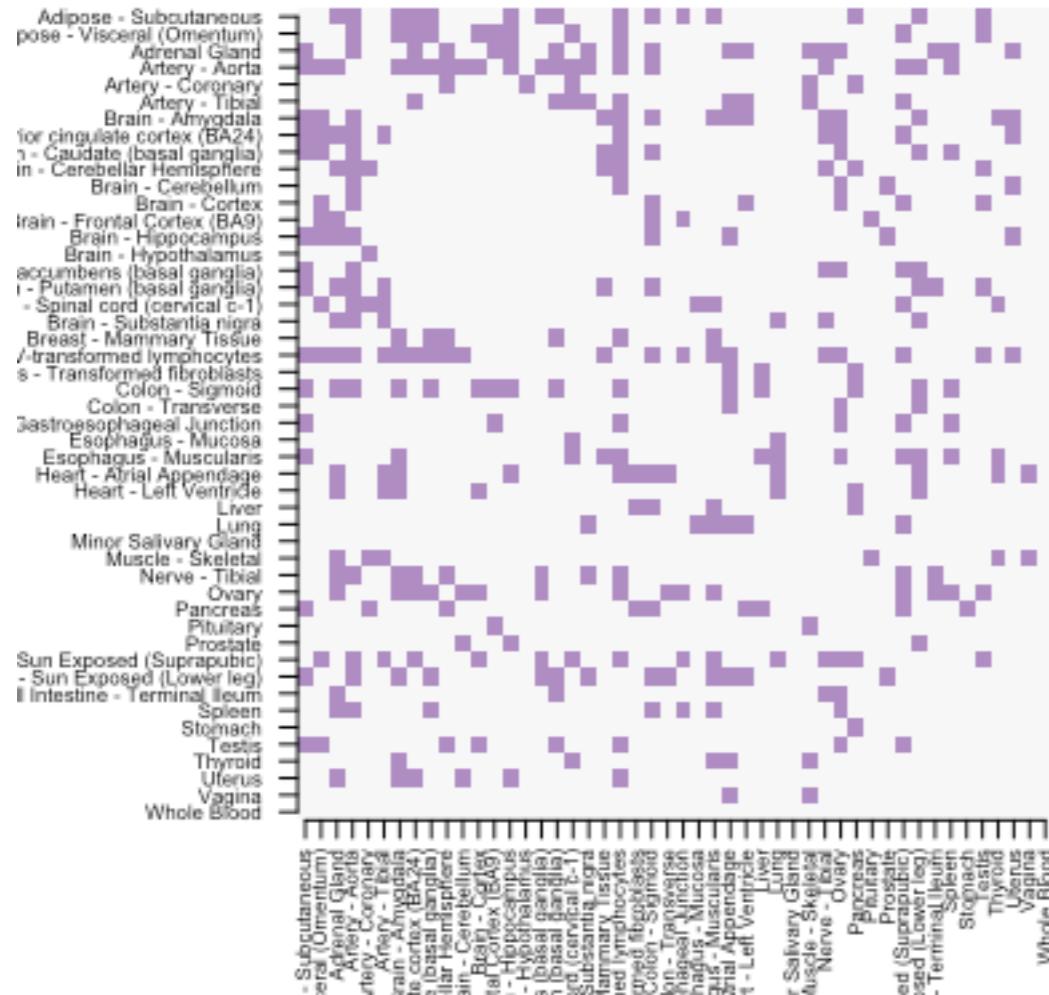


Which method separates tissues better?

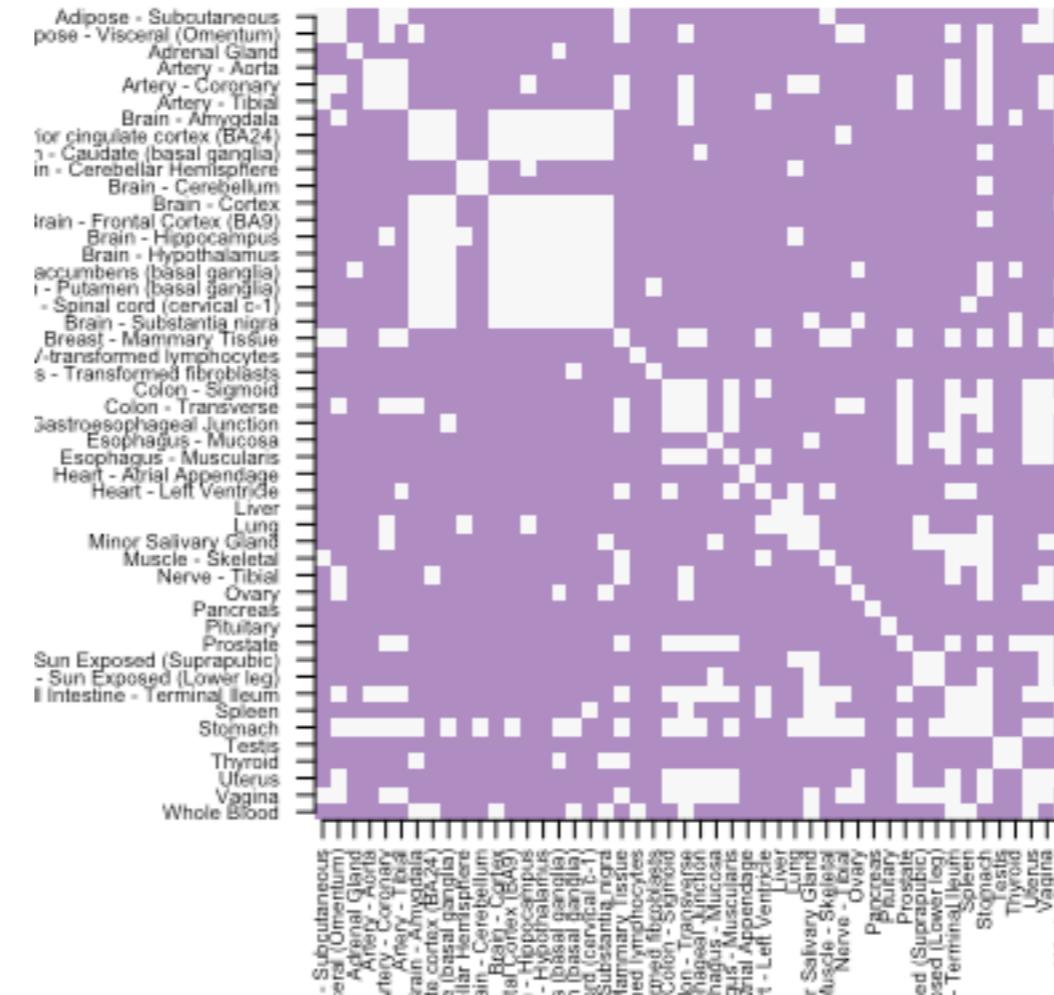
separates

does not separate

Hierarchical



Grade of Membership (GoM)

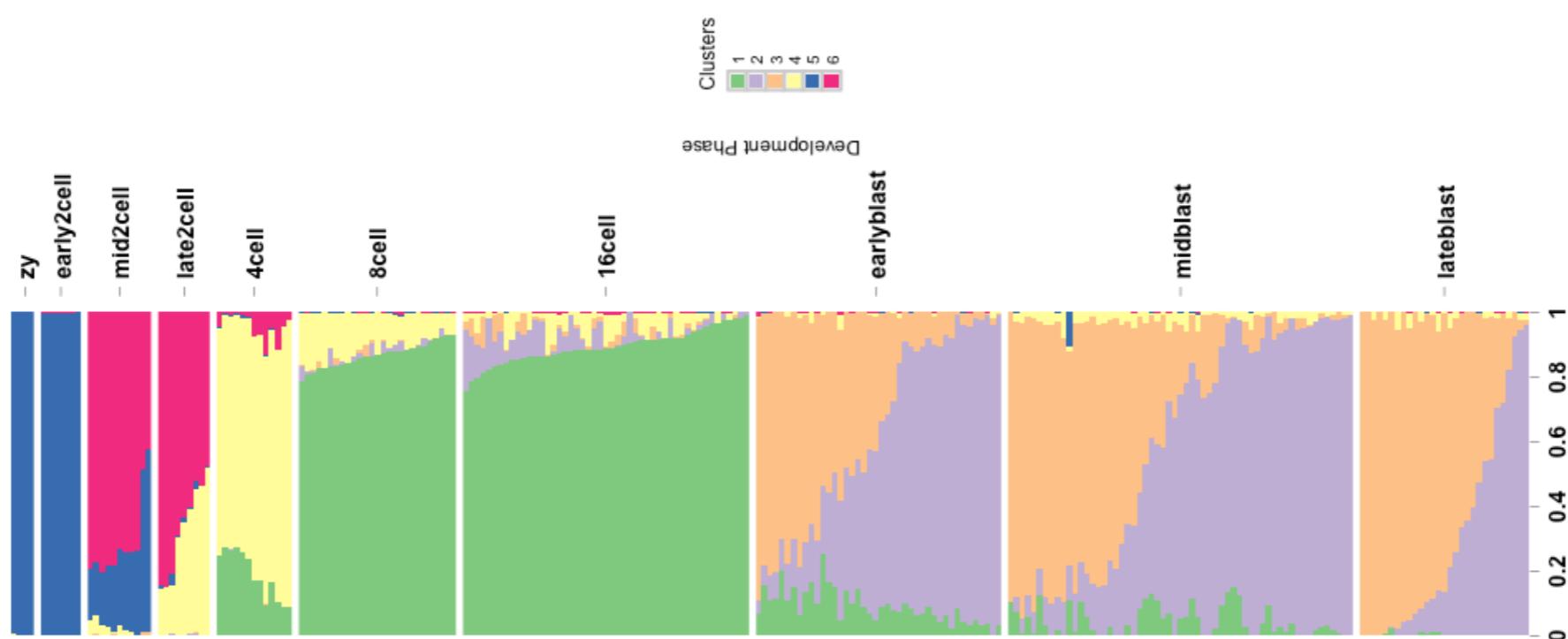
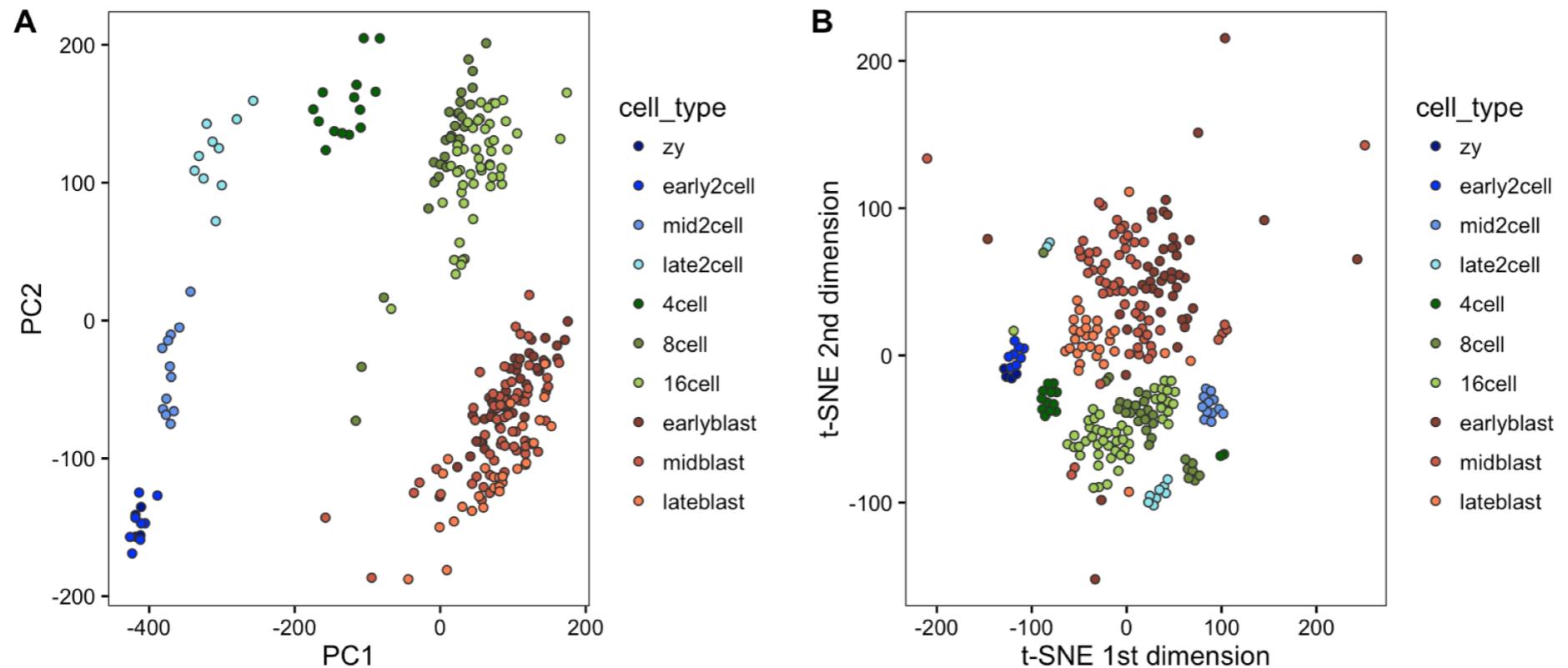


Single cell Example

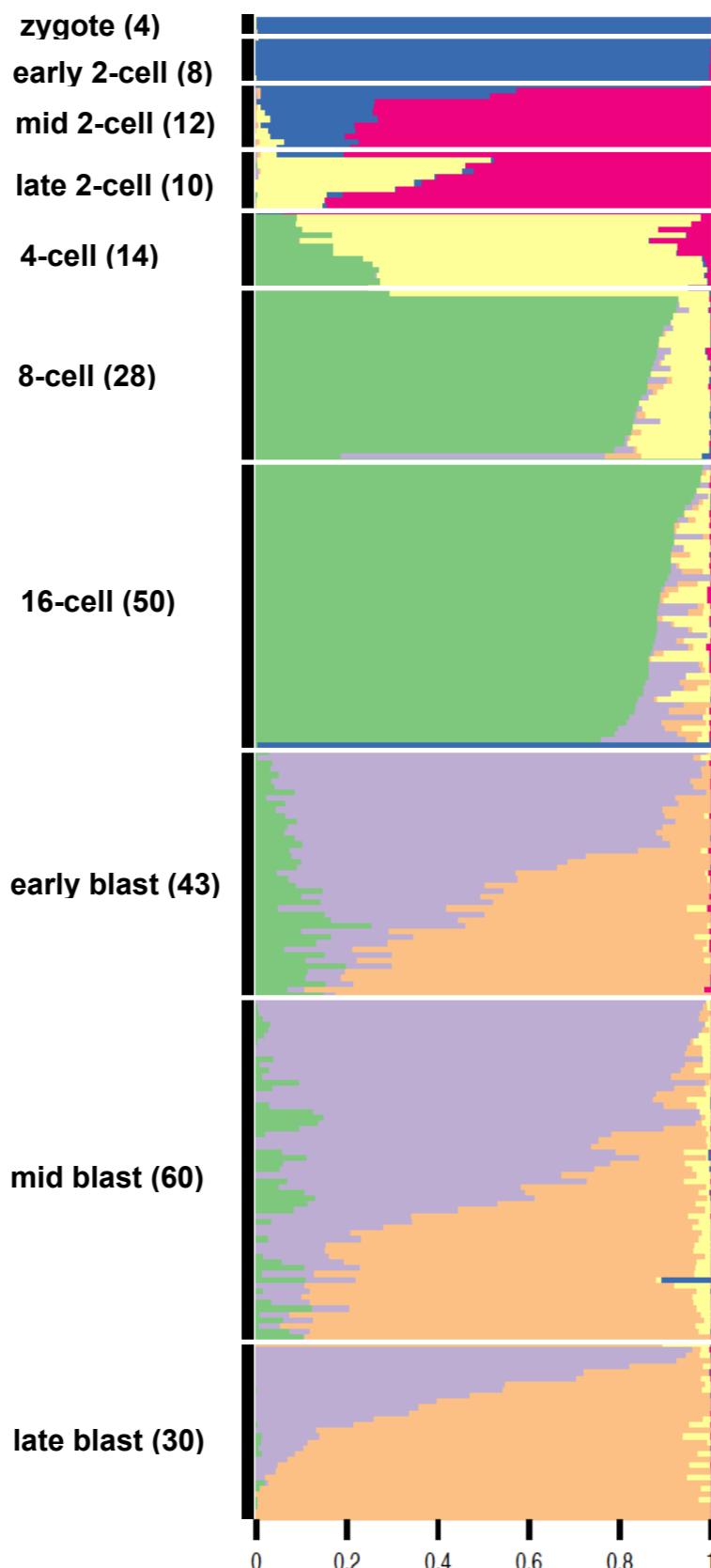
Deng et al (2014) collected single cell data across different developmental stages from zygote, 2-cell, 4-cell, 8-cell, 16-cell, and three stages of blastocysts (early, mid and late)

We plot admixture plot along with other dimension reduction methods like PCA and t-SNE

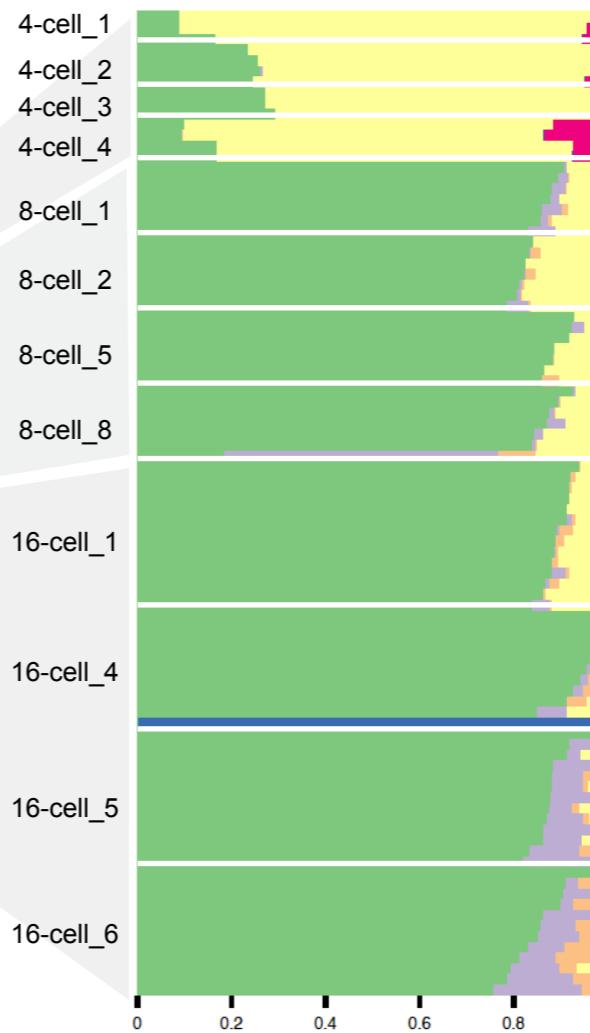
Deng Q1, Ramsköld D, Reinius B, Sandberg R. *Single-cell RNA-seq reveals dynamic, random monoallelic gene expression in mammalian cells.* Science. 2014 Jan 10;343(6167):193-6



STRUCTURE by developmental phase



STRUCTURE by embryo



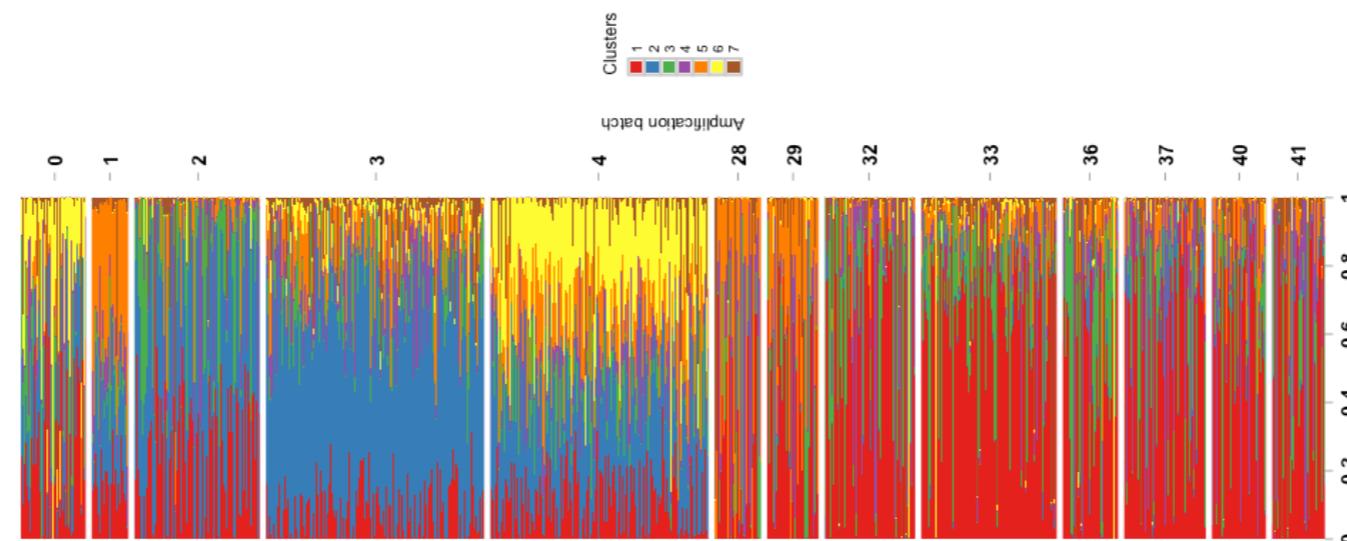
Clusters	Top 5 driving genes	GO annotations summary
1: Blue	<i>Bcl2l10, Tcf1, E330034G19Rik, LOC100502936, Oas1d</i>	<i>Gamete generation, Oocyte development & differentiation, Cell devision regulation</i>
2: Magenta	<i>Obox3, Zfp352, Gm8300, Usp17l5, BB287469</i>	<i>Nuclear body, Centriole, Microtubule organizing center part</i>
3: Yellow	<i>Rtn2, Ebna1bp2, Zfp259, Nasp, Cenpe</i>	<i>Nuclear part, Nuclear lumen, Intracellular organelle, Nucleolus, Microtubule</i>
4: Green	<i>Timd2, Isyna1, Alpl12, Prame15, Hsp90Ab1</i>	<i>Cytosol, Cytoplasmic Part, Catabolic Process, Protein Binding, Ubiquitin Ligase</i>
5: Purple	<i>Upp1, Tdgf1, Aqp8, Fabp5, Tat</i>	<i>Extracellular Exosome/Vesicle, Metabolic Process, Blood Vessel Morphogenesis</i>
6: Orange	<i>Actb, Krt18, Fabp3, Id2, Tspan8</i>	<i>Extracellular Exome/Organelle, Membrane Bound Organelle, Actin Cytoskeleton, Actin Filament Bundle</i>

Discussion

The results from the clustering method may be driven by partly by batch effects, besides the biological effects.

For data where the batches are confounded with biological subgroups, it becomes difficult to separate the two.

We encountered this issue in the analysis of *Jaitin et al (2014)* scRNA-seq data from mouse spleen.



Discussion

The method presented is unsupervised.

However, many recent studies sequence single cells from different cell types. These single cell information can be used to drive the clustering for bulk-RNA data.

In such a case, the clusters obtained would conform more with the cell-types by design and that was our primary interest.

Manuscript: **Clustering RNA-seq expression data using grade of membership models** (<http://biorxiv.org/content/early/2016/05/03/051631>)

R Package: *CountClust*

Bioconductor release version

(<https://www.bioconductor.org/packages/3.3/bioc/html/CountClust.html>)

Under development version

(<https://github.com/kkdey/CountClust>)

Acknowledgements

Matt Taddy (UChicago Booth School),
Jonathan Pritchard (Stanford University),
Yoav Gilad (UChicago Human Genetics)
Amos Tanay, Effi Kenigsberg,
Po Yuan Tung, John Blischak (Gilad Lab)
Gao Wang, Raman Shah, Nan Xiao (Stephens Lab)

Co-authors



Joyce Hsiao



Matthew Stephens