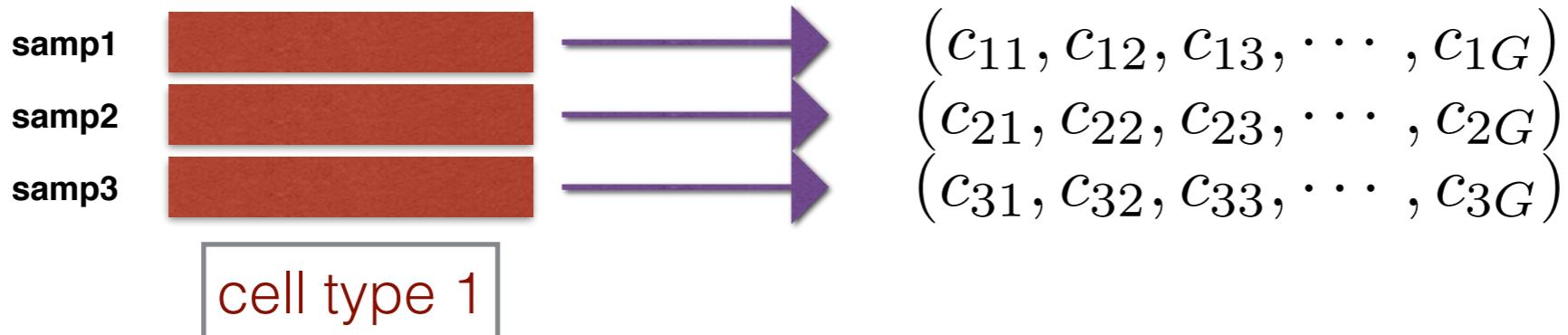
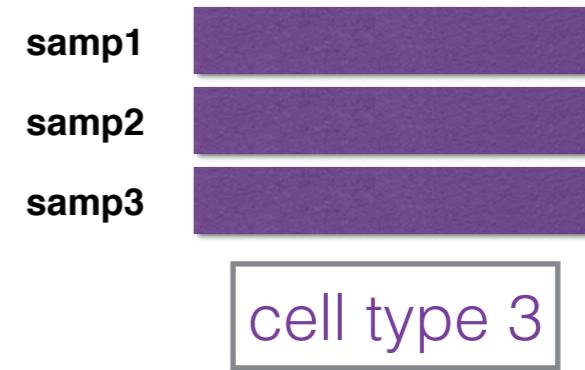
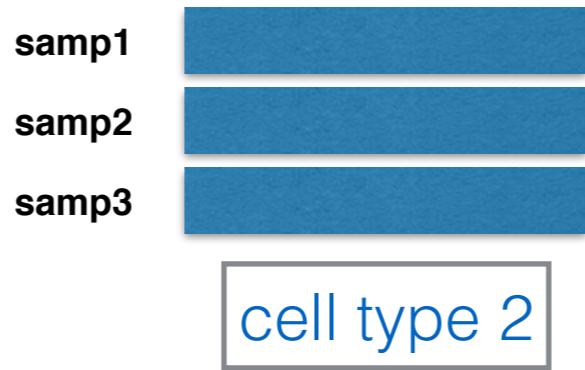
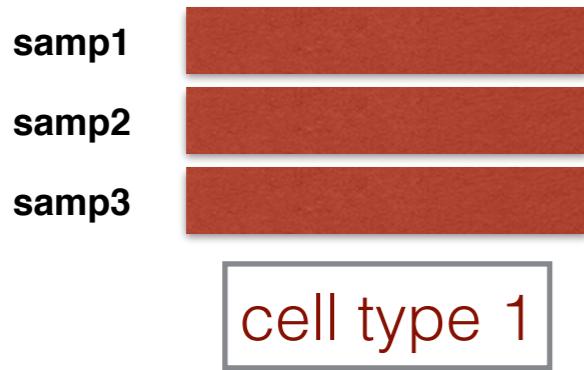


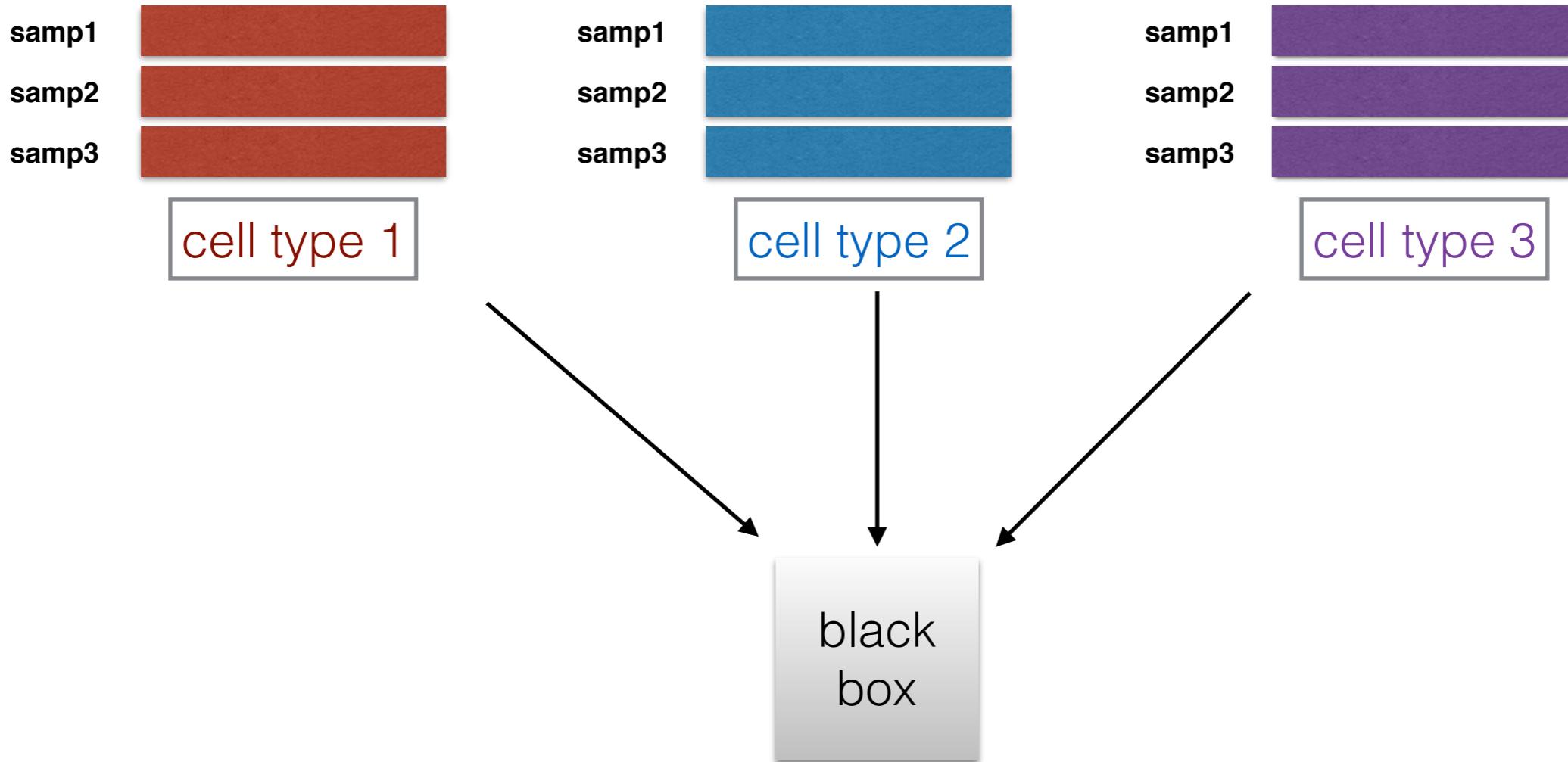
# **Grade of membership classification in RNA-seq data**

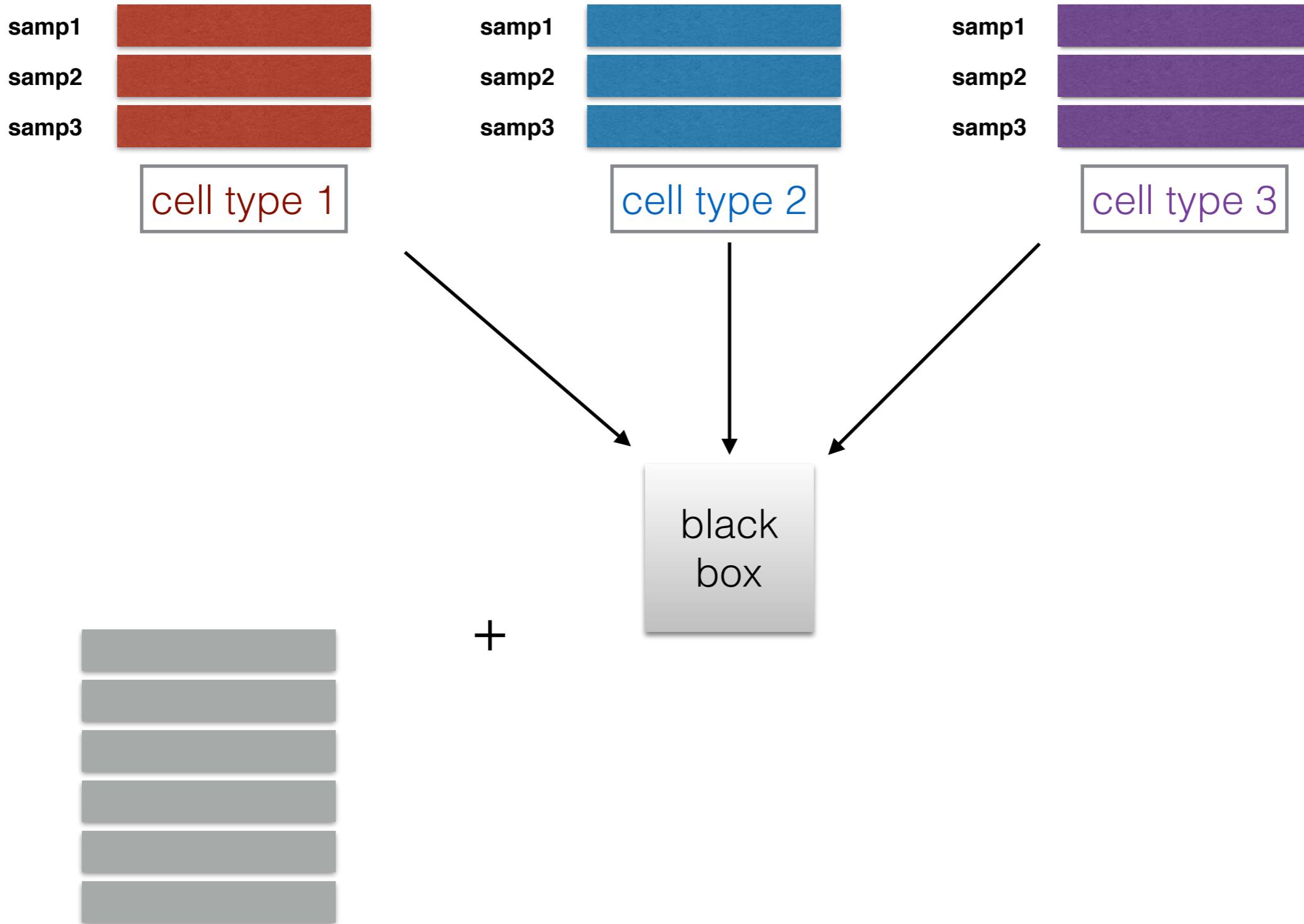
- Model intuition
- Statistical framework
- Workflow for RNA-seq data
- GTEx applications
- Single-cell applications
- Summary

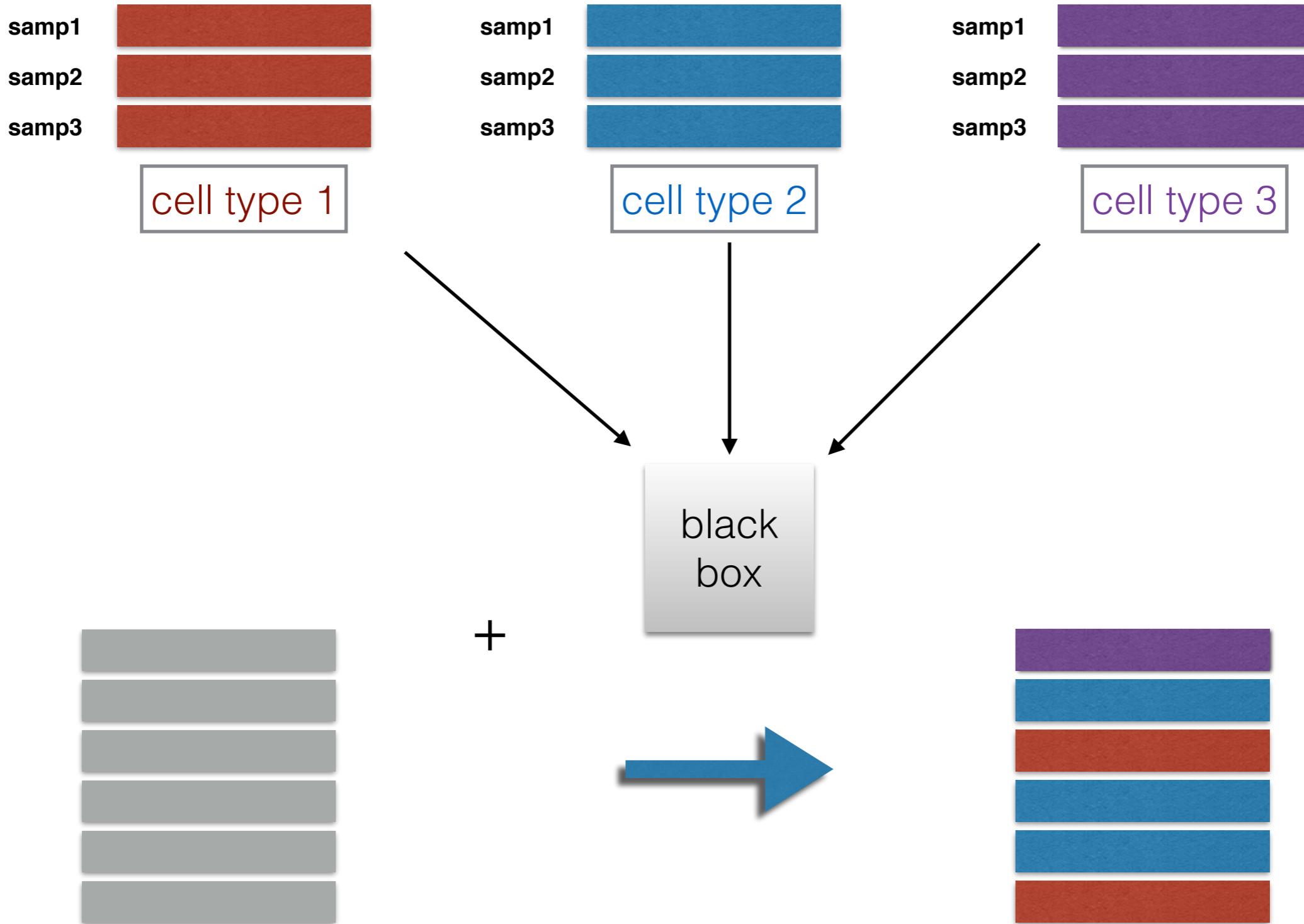


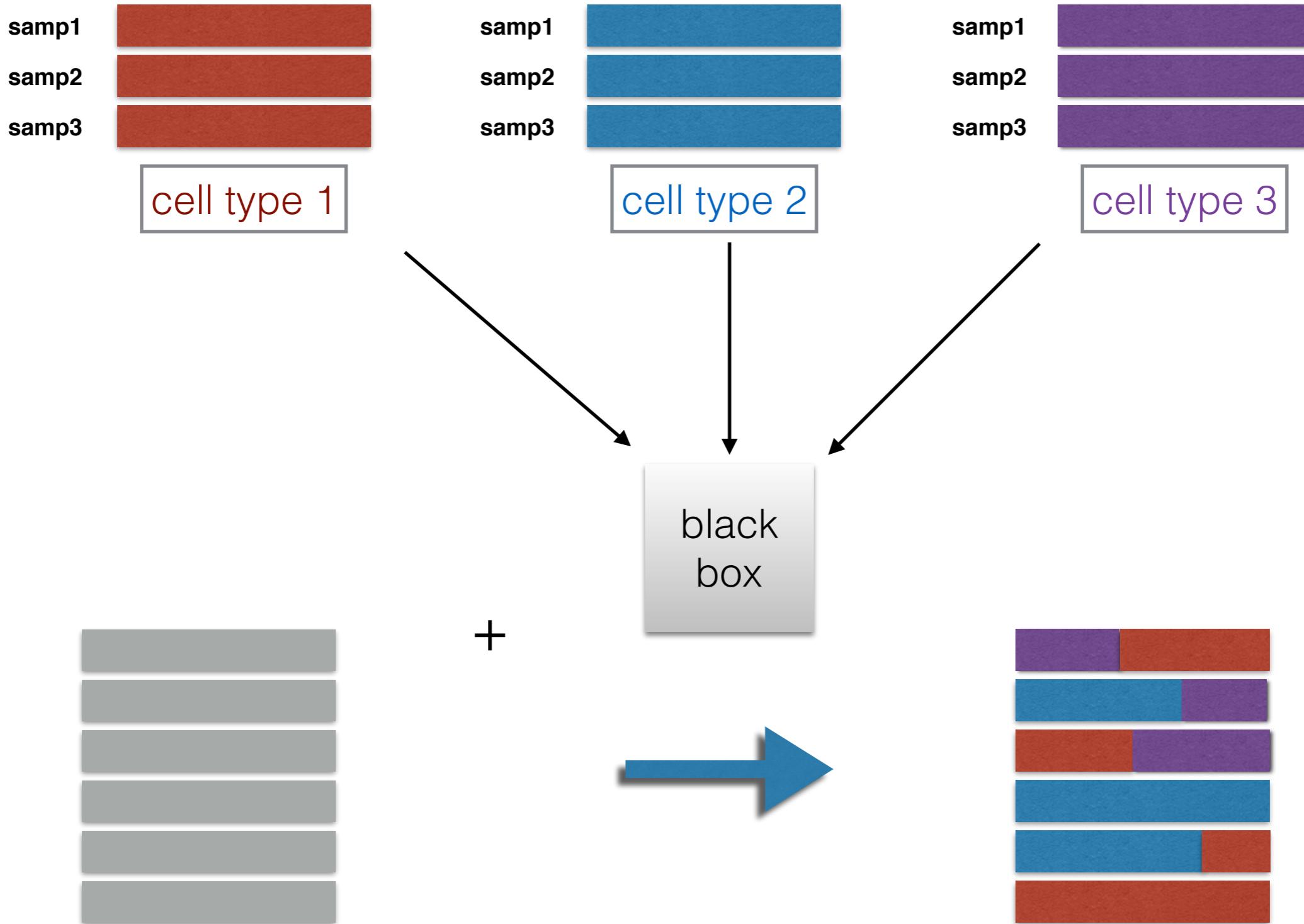
**RNA-seq reads count data  
across G genes for each  
sample**









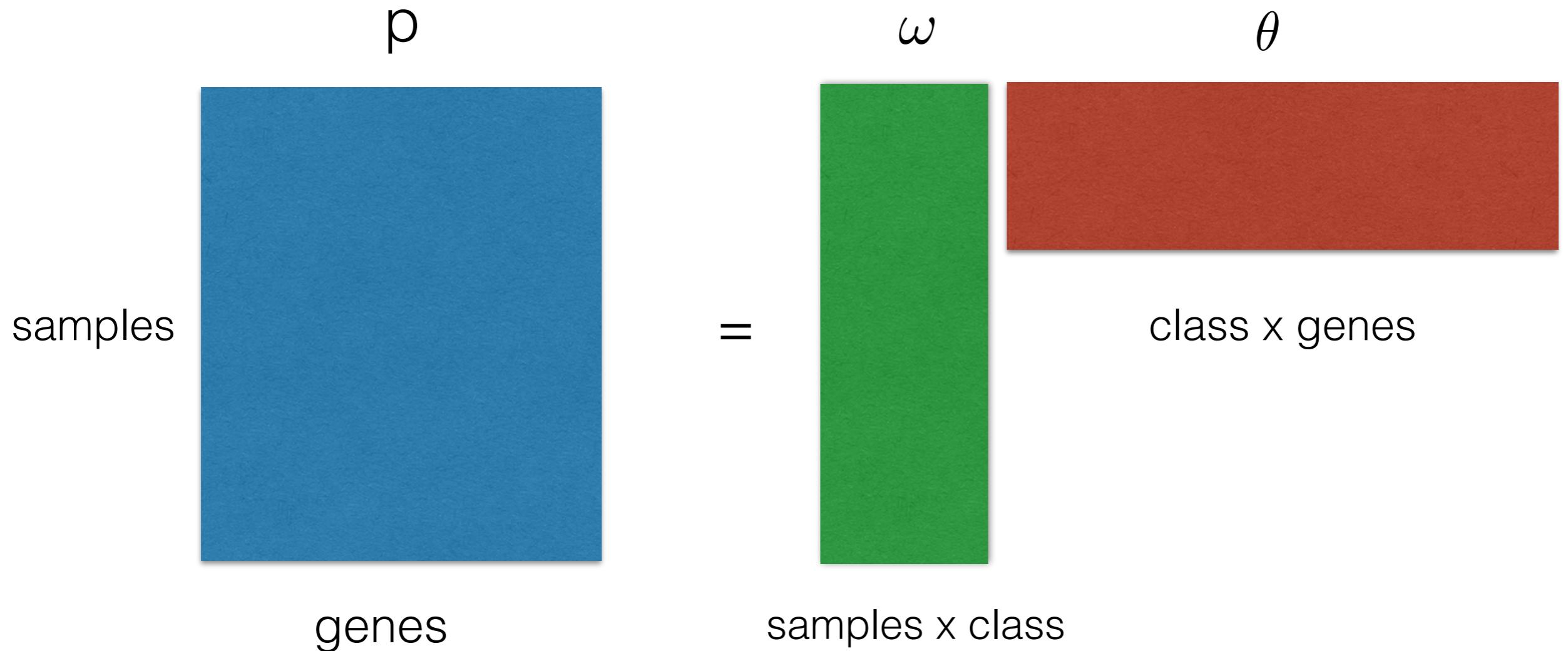


Grade of membership model

# The model

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

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



# **Unsupervised approach**

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

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

$\omega_{nk}$  grade of membership of sample n in cluster or topic k

$\theta_{kg}$  relative proportion of expression of gene g in cluster k

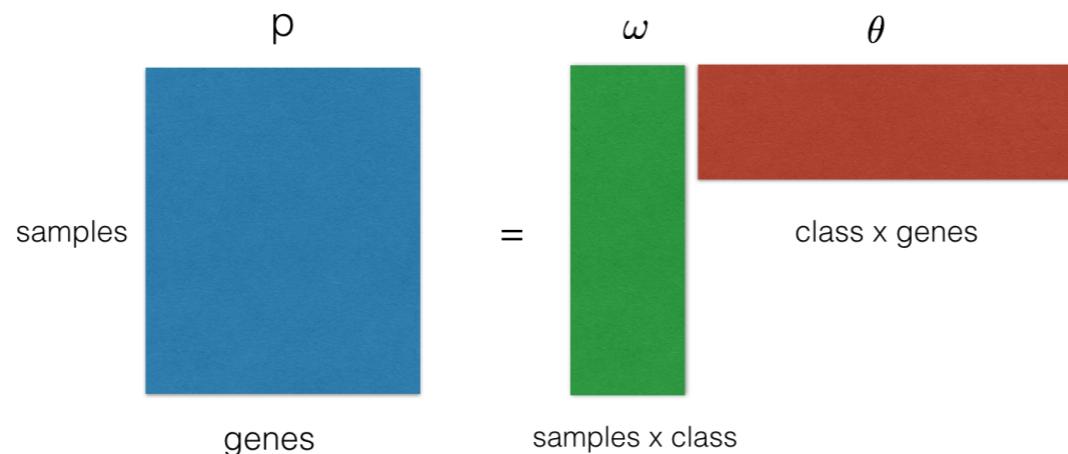
$\omega_{nk}$  and  $\theta_{kg}$  were determined in an unsupervised manner  
assuming Dirichlet priors

Note: R package **maptpx** was used to perform unsupervised GoM fitting.

# **classtpx: an semi-supervised approach**

In **omega.fix**, some samples coming from each cell type  $k$  are treated as training and we assign grade of membership of those samples in class  $k$  to be 1 and that in other classes to be 0.

In **theta.fix** method, we fix  $\theta_k$ . for cluster  $k$  by pooling the relative expression patterns across genes for training samples coming from cluster  $k$



The test data would be actual tissue samples which would be a mixture of these different cell types or classes.



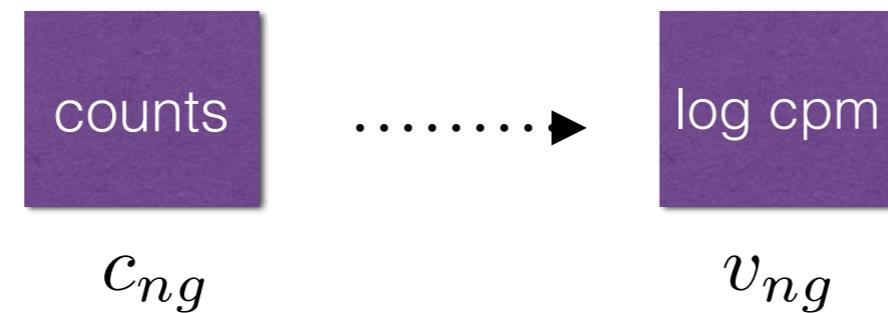
RNA-seq application

# Challenges

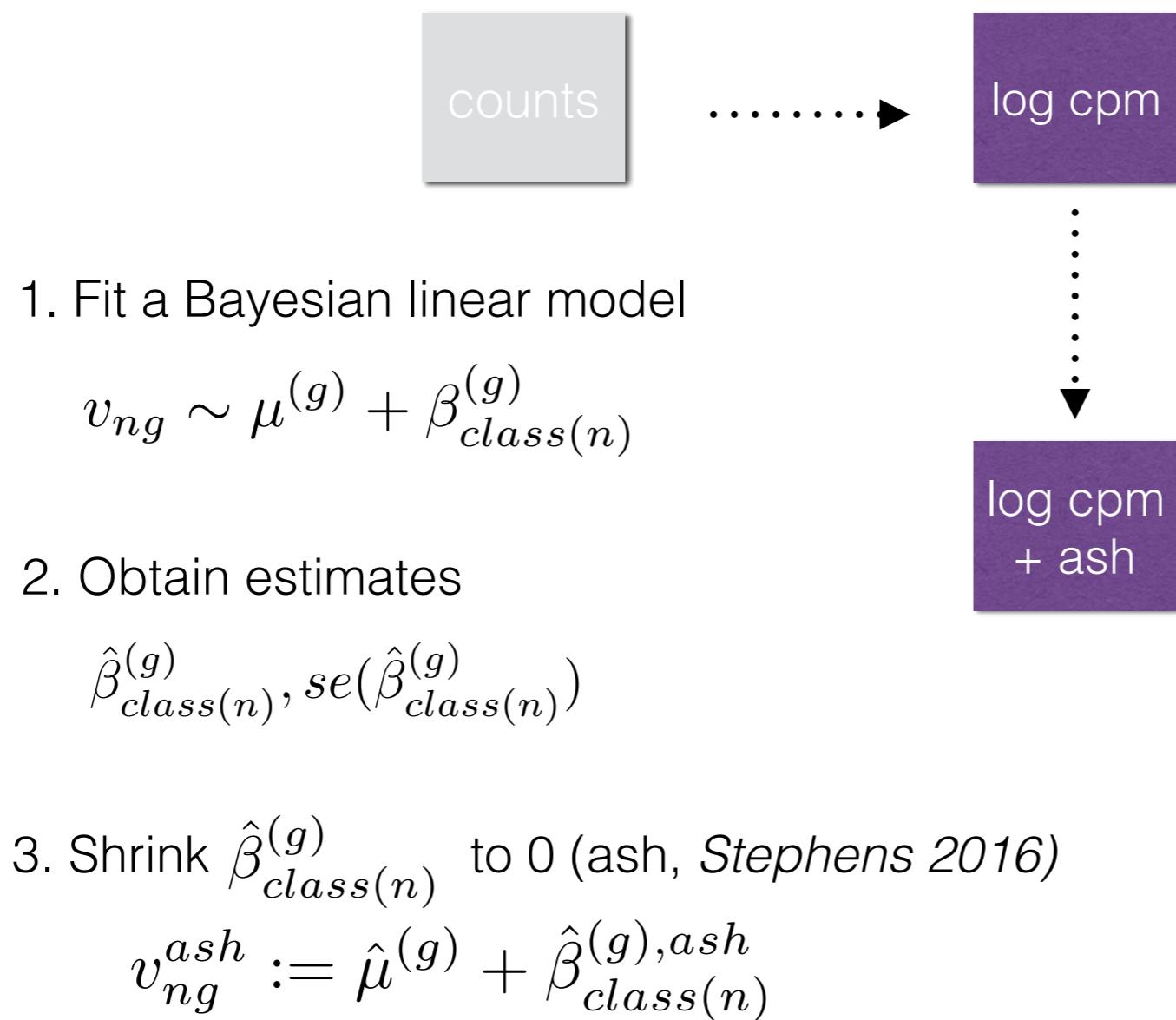
- Normalization
- Expression noise (should we shrink?? adaptively?)
- Count-based or log-CPM based analysis?

# **Workflow**

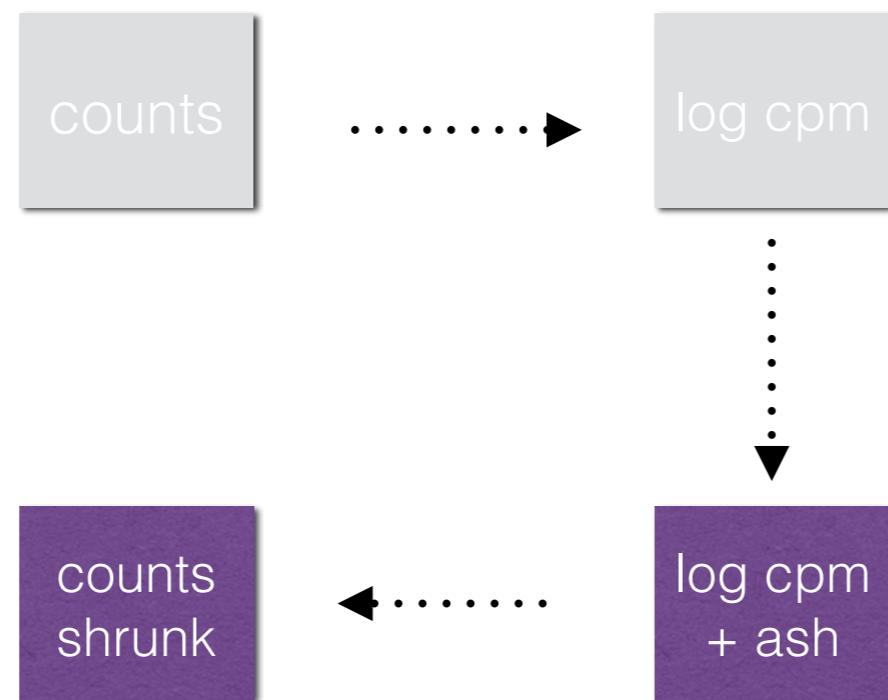
$$c_{ng} \rightarrow \log_2 \left( \frac{c_{ng} + 0.5}{c_{n+} + 1} \times 10^6 \right)$$



# Workflow



# Workflow



4. Reverse transform

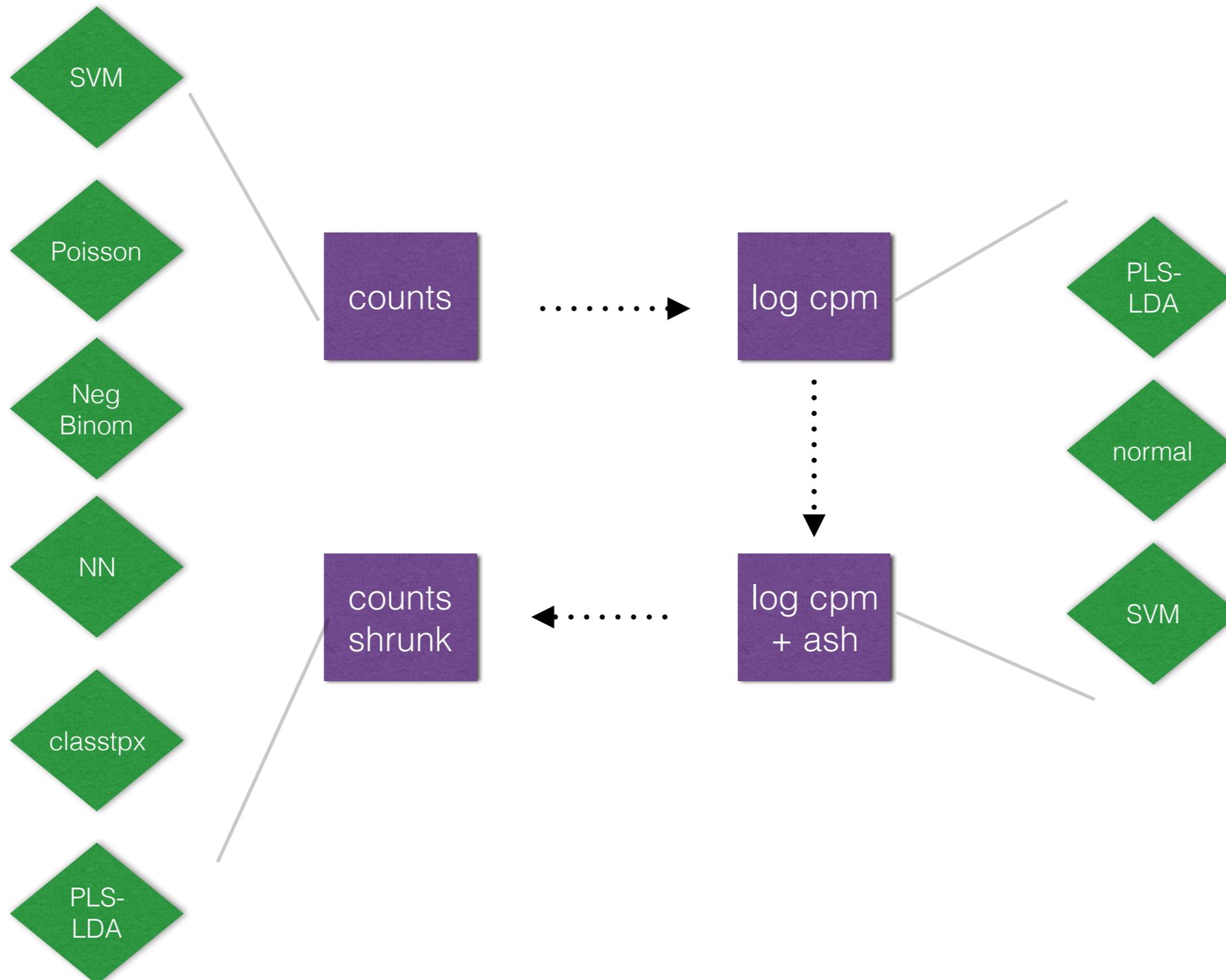
$$v_{ng}^{ash} \xrightarrow{rev.voom} c_{ng}^{ash}$$

# **Classification in RNA-seq data**

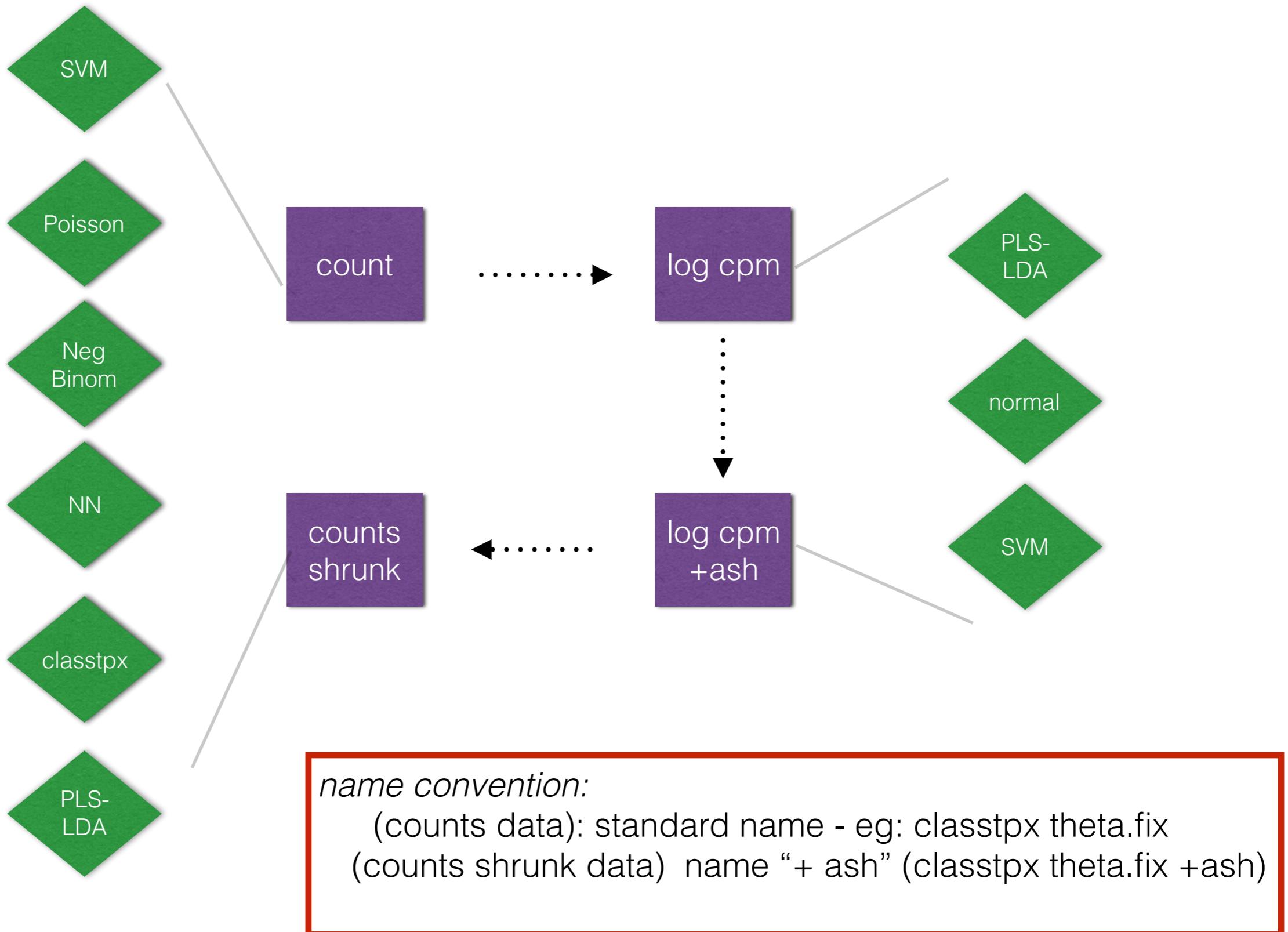
1. Model-based versus non-model based
2. Soft versus hard classification

- Non-model based approach:  
*Support Vector machines, PLS-LDA, Neural network*
- Model based approach  
*Poisson model on counts, Negative Binomial model on counts, Normal model, **grade of membership approach***

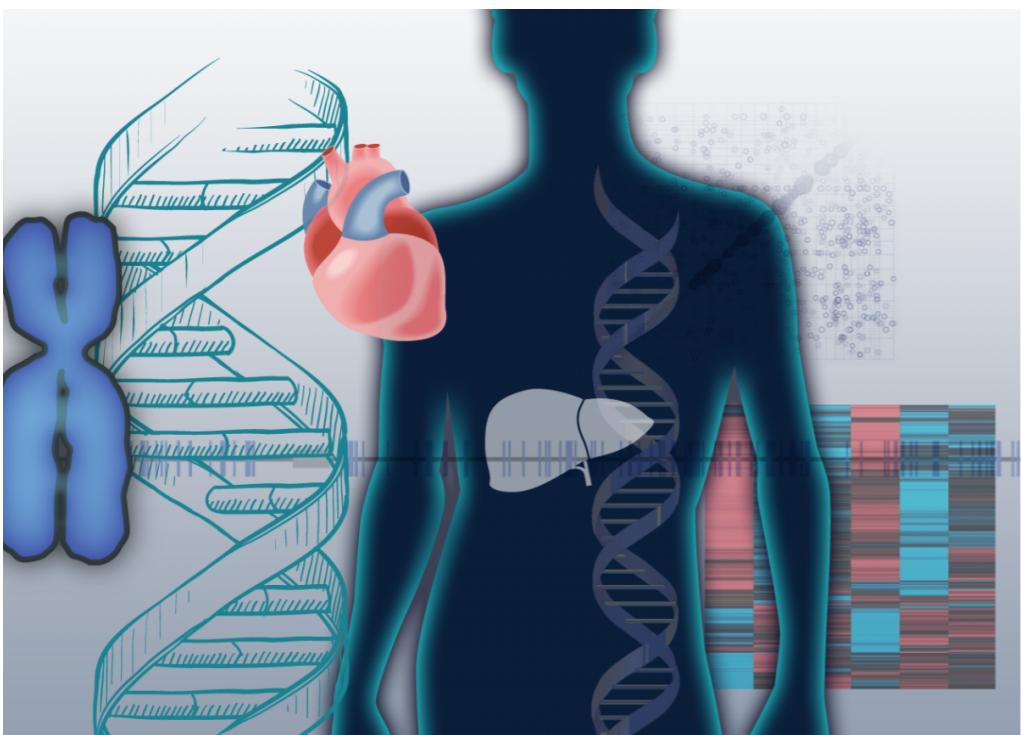
# Workflow



# Workflow



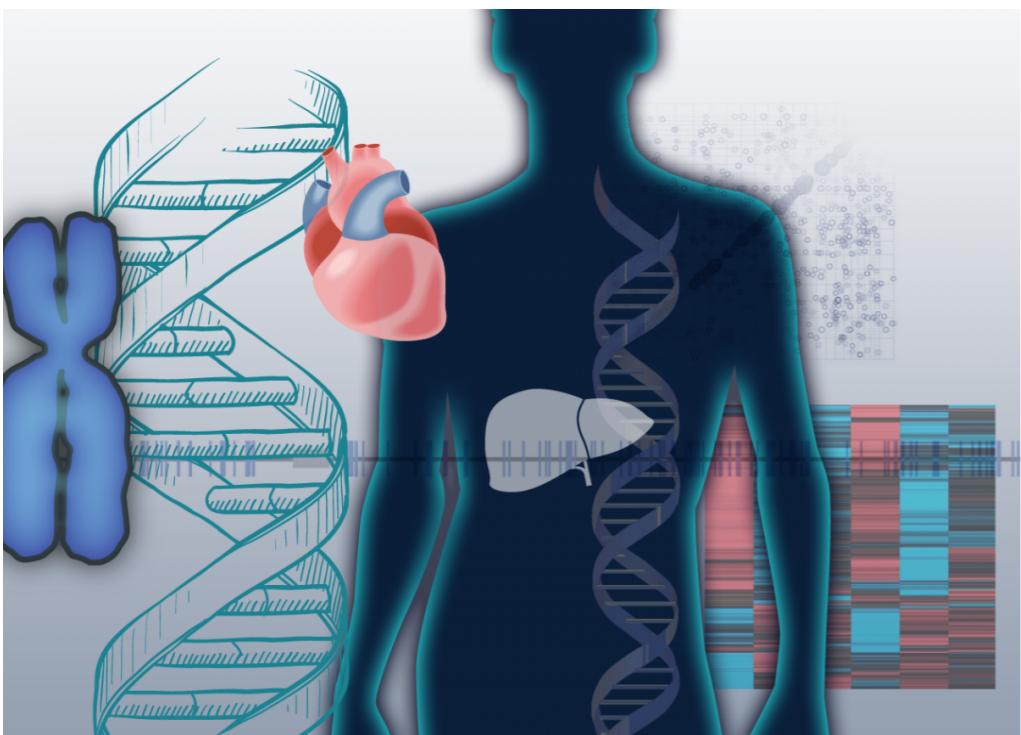
# Results Summary - GTEx bulk RNA-seq



# Genotype Tissue Expression Project (GTEx V6)

tissue sample RNA-seq data from 53 tissues (total 8555 samples).

We filtered out 16,000 genes initially.



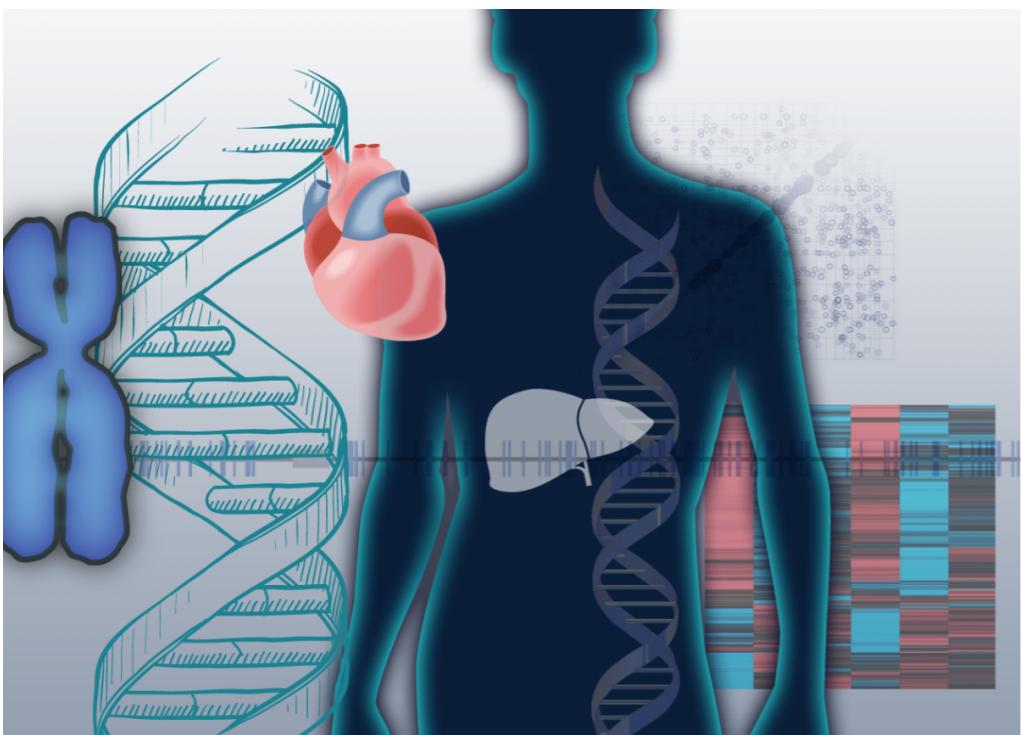
# Genotype Tissue Expression Project (GTEx V6)

tissue sample RNA-seq data from 53 tissues (total 8555 samples).

We filtered out 16,000 genes initially.

For classtpx, we selected those tissues with more than 100 samples- there were 36 such tissues.

50 samples per tissue were used as training.



# Genotype Tissue Expression Project (GTEx V6)

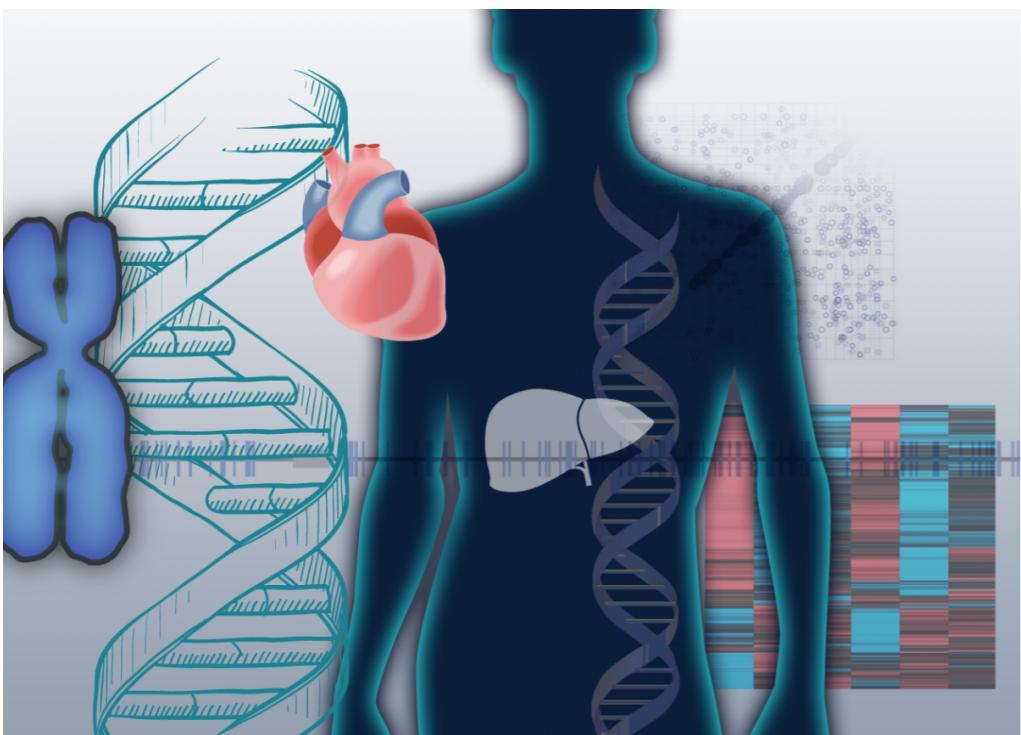
tissue sample RNA-seq data from 53 tissues (total 8555 samples).

We filtered out 16,000 genes initially.

For classtpx, we selected those tissues with more than 100 samples- there were 36 such tissues.

50 samples per tissue were used as training.

For classification on all 36 tissues, we extracted 7000 top genes distinguishing the 36 tissue classes using *CountClust::ExtractTopfeatures()* [except neural network for which top 2000 genes were used]



# Genotype Tissue Expression Project (GTEx V6)

tissue sample RNA-seq data from 53 tissues (total 8555 samples).

We filtered out 16,000 genes initially.

For classtpx, we selected those tissues with more than 100 samples- there were 36 such tissues.

50 samples per tissue were used as training.

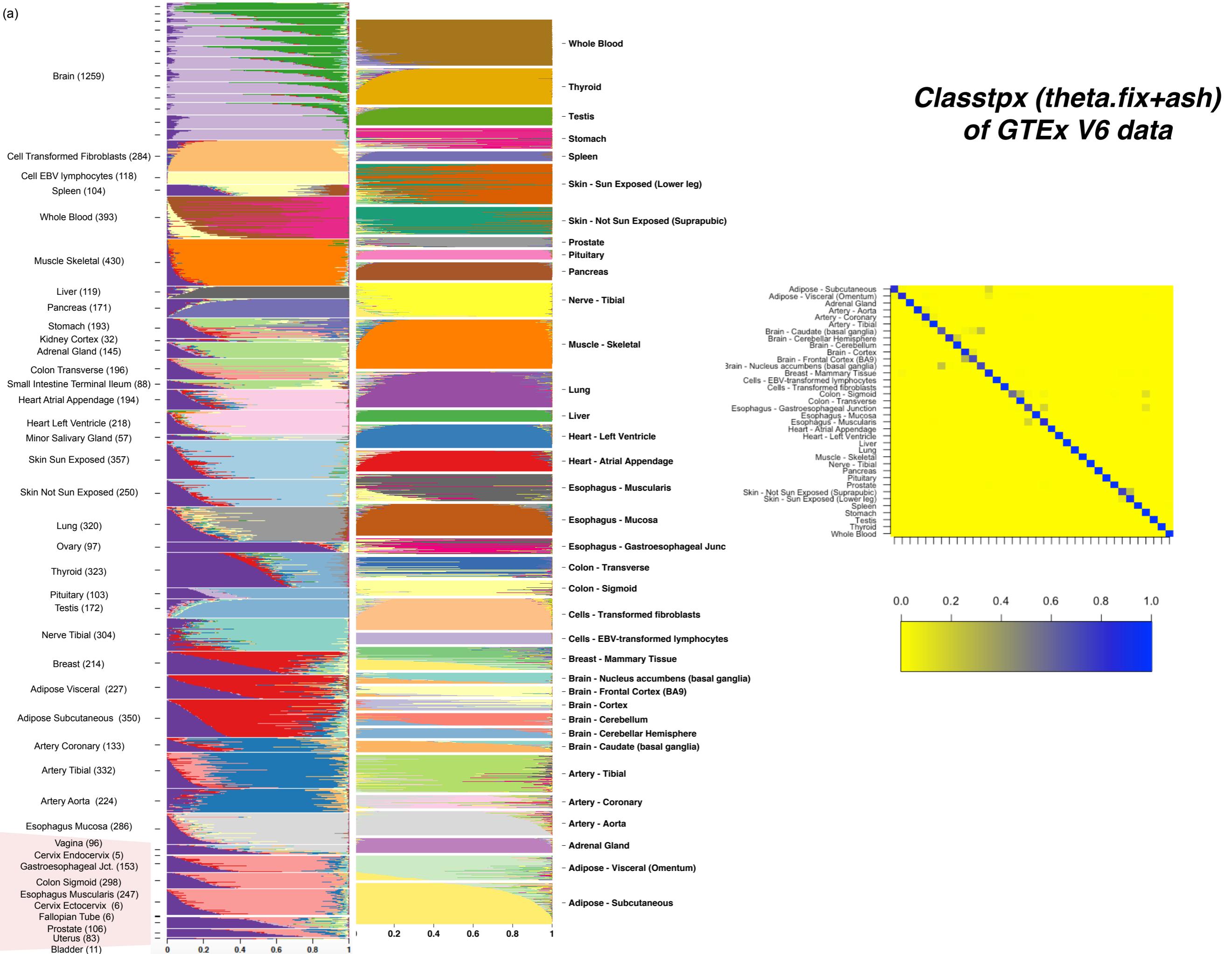
For classification on all 36 tissues, we extracted 7000 top genes distinguishing the 36 tissue classes using *CountClust::ExtractTopfeatures()* [except neural for which top 2000 genes were used]

Separate classification was done on pairs and triplets of tissues

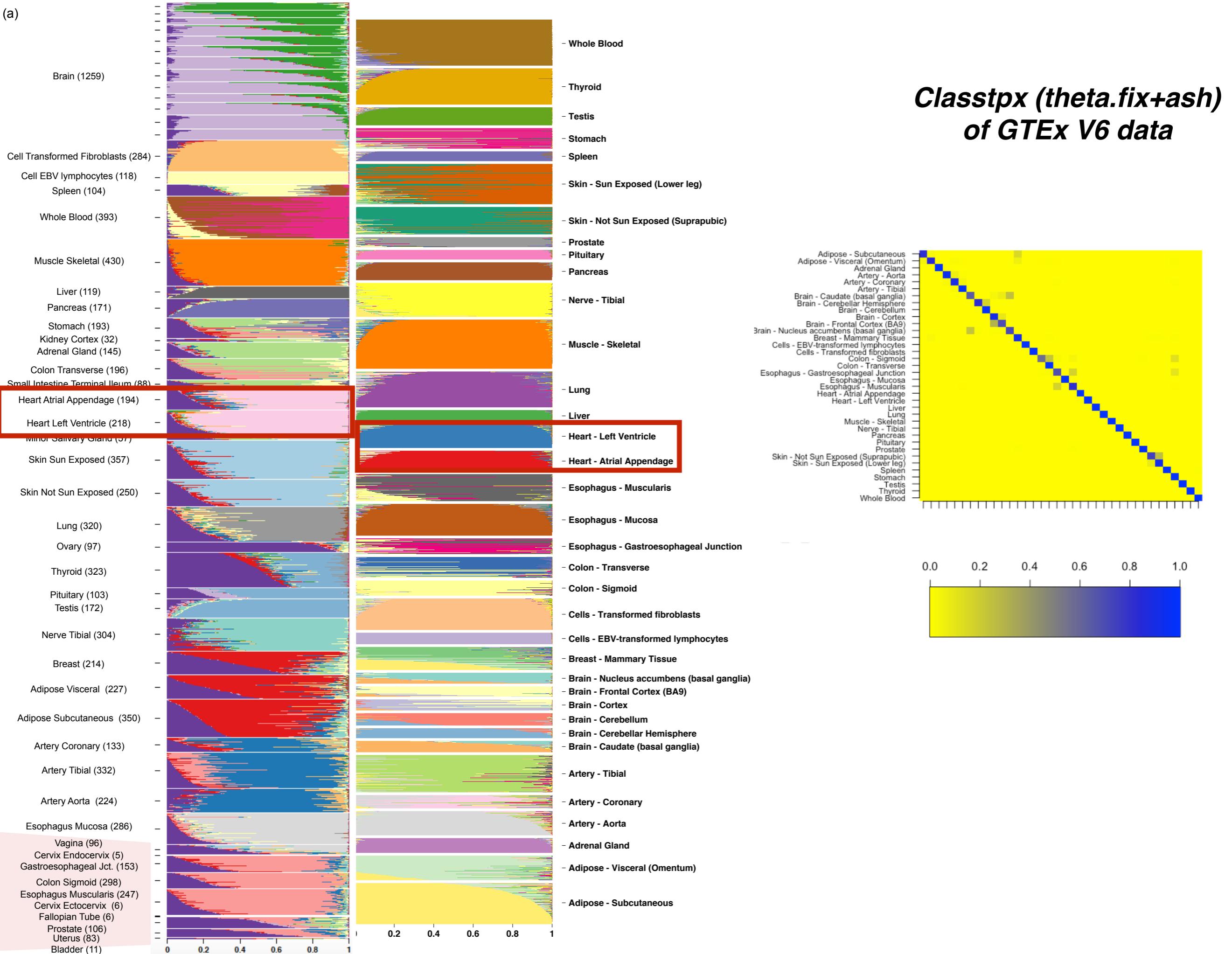
- Arteries (Coronary, Tibial, Aorta),
- Breast Mammary Tissue, Adipose Subcutaneous, Adipose Visceral

For these comparison studies we did analysis on all genes and on top 3000 genes

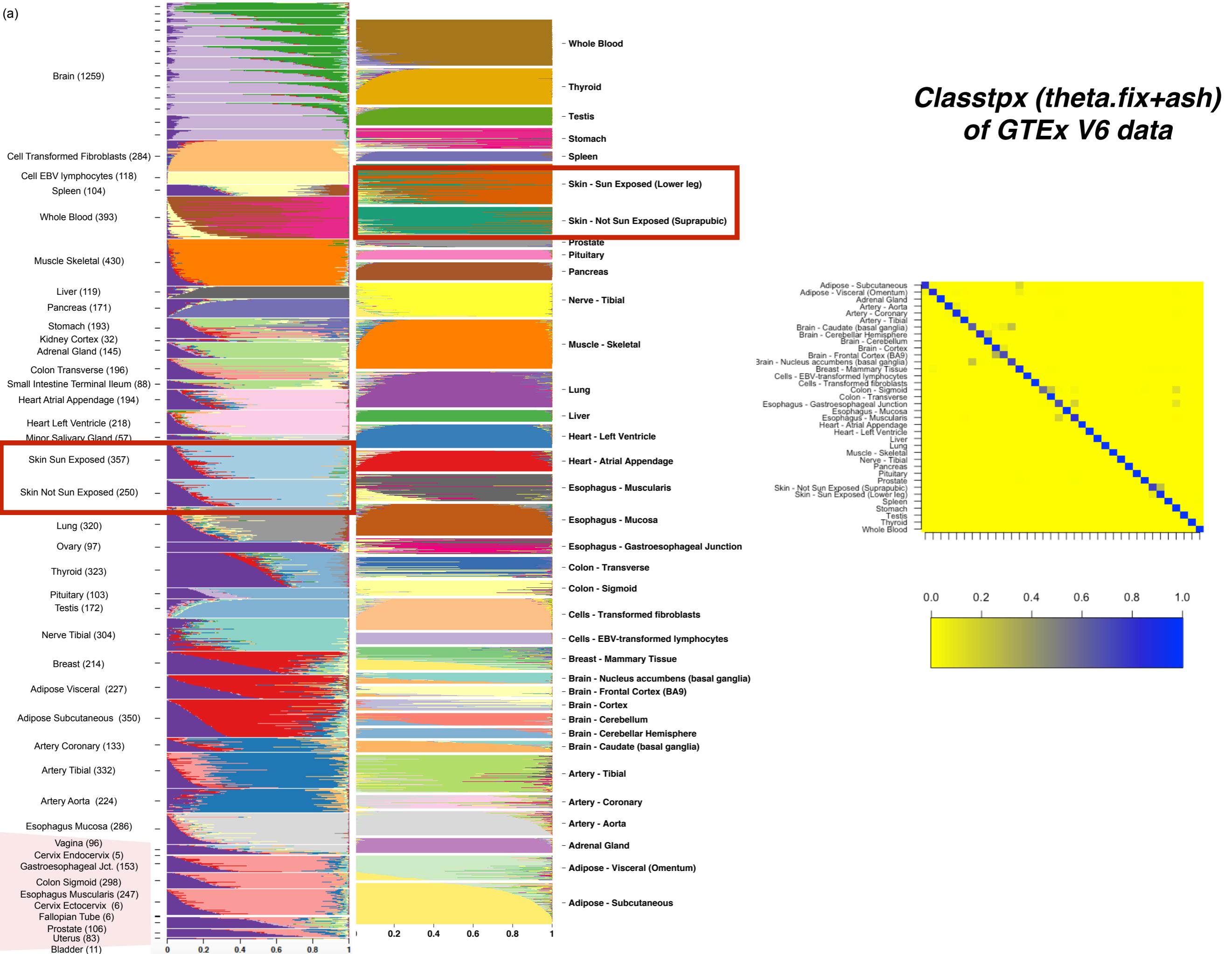
(a)



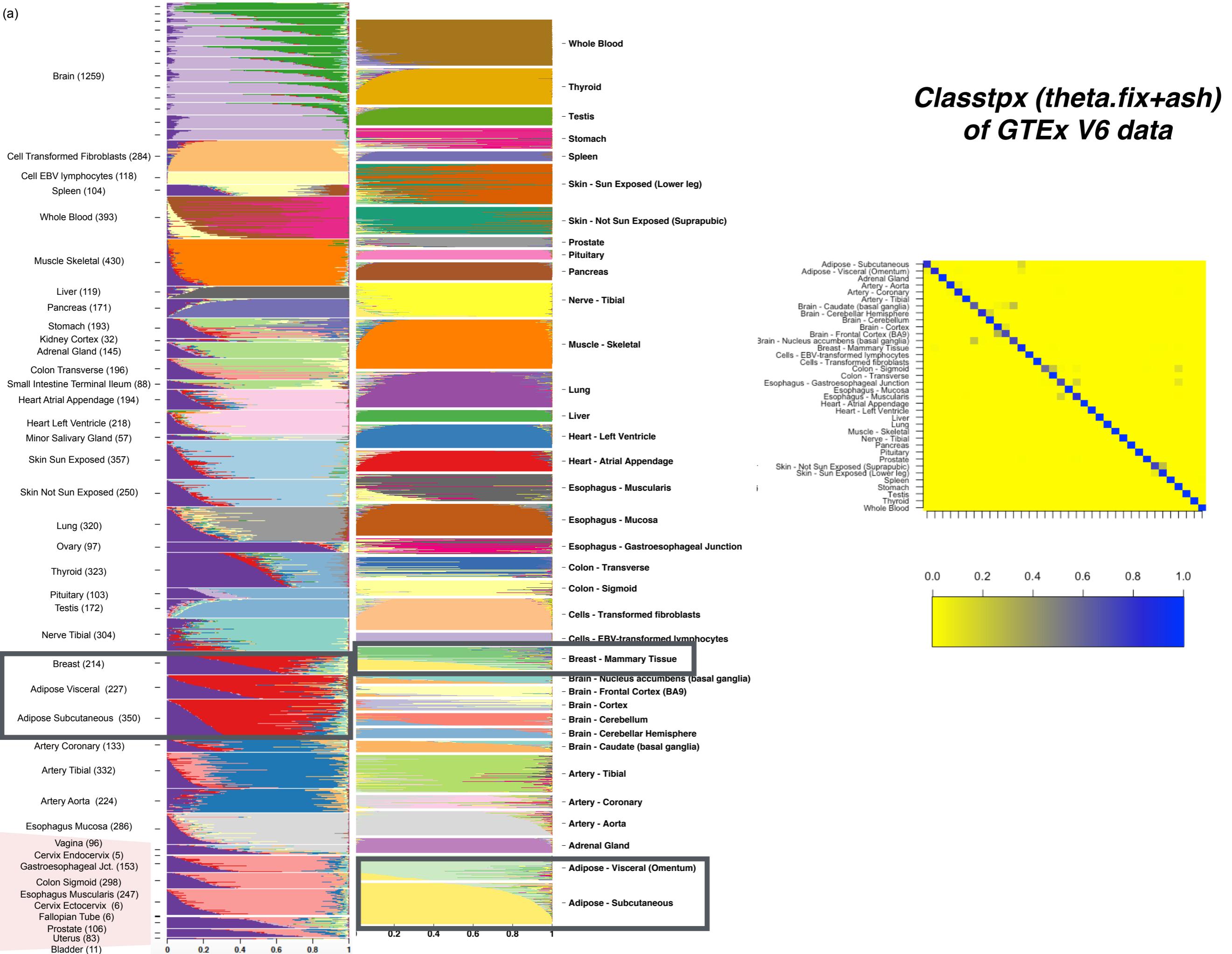
(a)



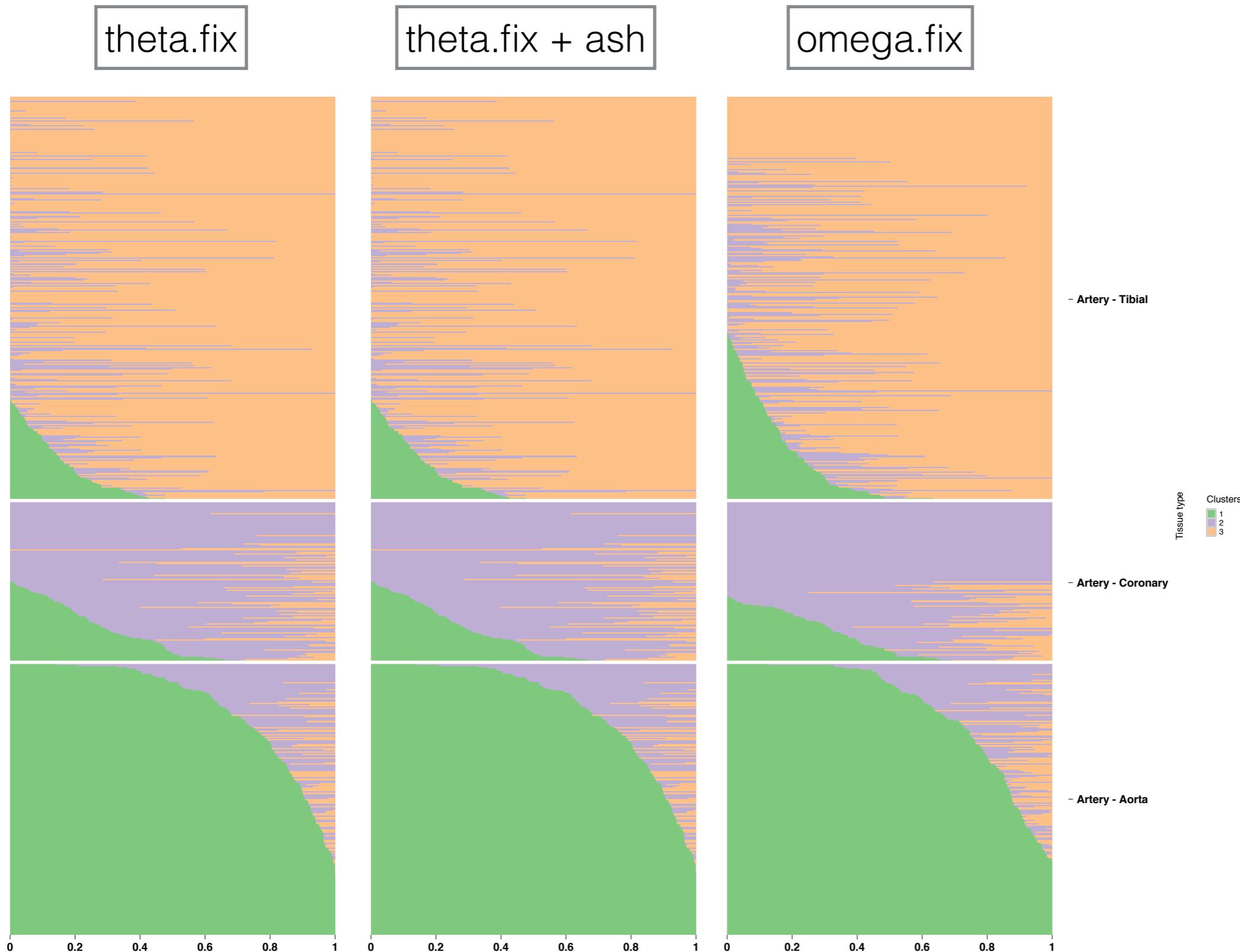
(a)



(a)



# *Classpx of three artery tissues in GTEx V6 data*



## ***classtpx vs Other classifiers***

To compare with hard clustering methods like SVM, PLS-LDA etc, we assign each sample to the cluster or class that is most representative of the sample (which class has the highest grade of membership for that sample)

misclass prop	All tissues	Arteries (3 tissues)	Arteries3000 (3 tissues)	Breast +Adipose	Breast + Adipose 3000
<b>classtpx theta.fix</b>	0.07	0.09	0.08	0.14	0.13
<b>classtpx theta.fix +ash</b>	0.07	0.09	0.08	0.13	0.13
<b>classtpx omega.fix</b>	0.09	0.13	0.11	0.22	0.19
<b>SVM (counts)</b>	0.09	0.22	0.20	0.21	0.15
<b>SVM (log cpm)</b>	0.05	0.05	0.03	0.08	0.05
<b>SVM (log cpm+ash)</b>	0.05	0.05	0.03	0.09	0.05
<b>PLS-LDA (counts)</b>	0.14	0.06	0.03	0.11	0.10
<b>PLS-LDA (log cpm)</b>	0.03	0.22	0.2	0.24	0.06
<b>PLS-LDA (log cpm +ash)</b>	0.03	0.16	0.16	0.18	0.06
<b>Normal</b>	0.09	0.12	0.09	0.18	0.11
<b>Poisson</b>	0.08	0.09	0.08	0.14	0.13
<b>Neg Binom</b>	0.10	0.18	0.10	0.23	0.17
<b>NN</b>	0.21	-	-	-	-

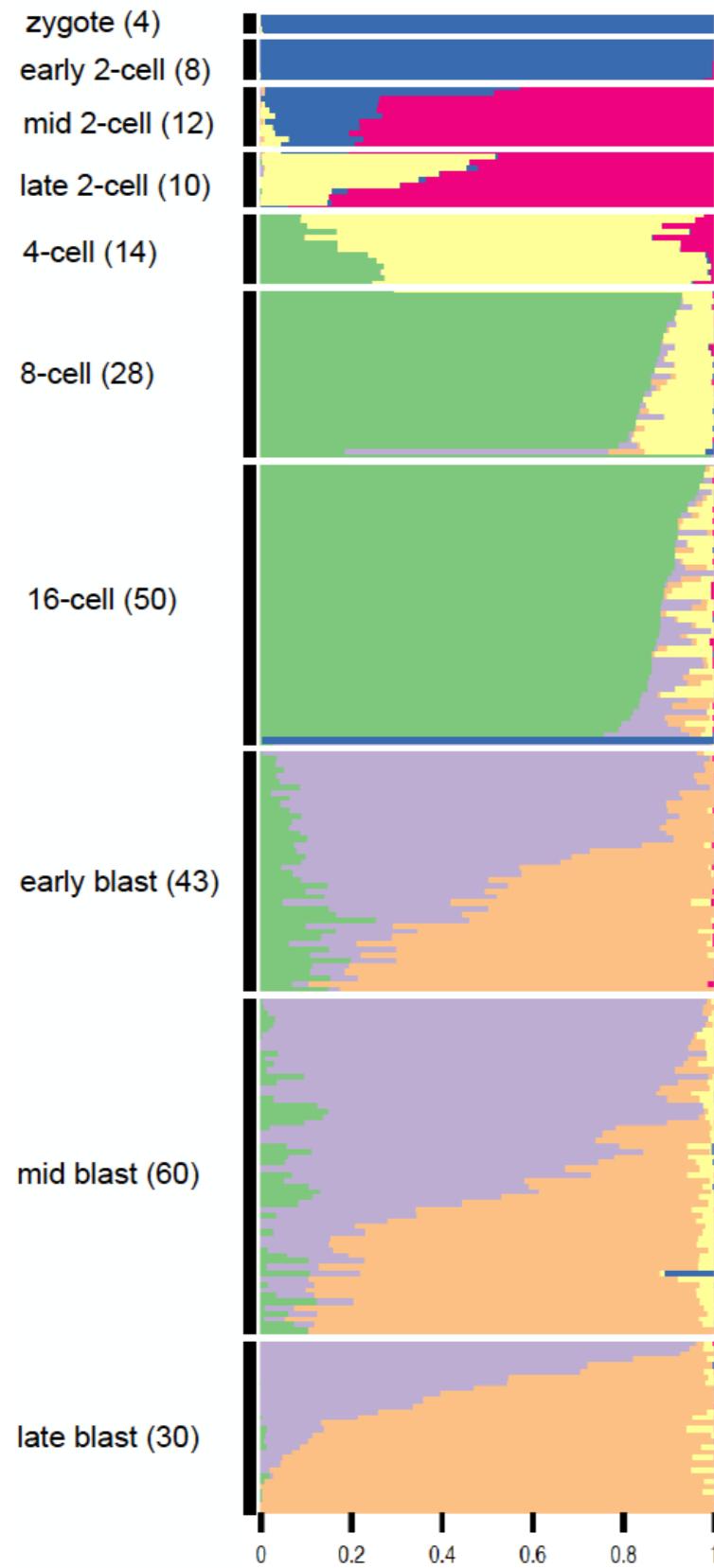
	All tissues	Arteries (3 tissues)	Arteries3000 (3 tissues)	Breast +Adipose	Breast + Adipose 3000
<b>classtpx theta.fix</b>	0.07	0.09	0.08	0.14	0.13
<b>classtpx theta.fix +ash</b>	0.07	0.09	0.08	0.13	0.13
<b>classtpx omega.fix</b>	0.09	0.13	0.11	0.22	0.19
<b>SVM (counts)</b>	0.09	0.22	0.20	0.21	0.15
<b>SVM (voom)</b>	0.05	0.05	0.03	0.08	0.05
<b>SVM (voom+ash)</b>	0.05	0.05	0.03	0.09	0.05
<b>PLS-LDA (counts)</b>	0.14	0.06	0.03	0.11	0.10
<b>PLS-LDA (voom)</b>	0.03	0.22	0.2	0.24	0.06
<b>PLS-LDA (voom +ash)</b>	0.03	0.16	0.16	0.18	0.06
<b>Normal</b>	0.09	0.12	0.09	0.18	0.11
<b>Poisson</b>	0.08	0.09	0.08	0.14	0.13
<b>Neg Binom</b>	0.10	0.18	0.10	0.23	0.17
<b>NN</b>	0.21	-	-	-	-

	All tissues	Arteries (3 tissues)	Arteries3000 (3 tissues)	Breast +Adipose	Breast + Adipose 3000
<b>classtpx theta.fix</b>	0.07	0.09	0.08	0.14	0.13
<b>classtpx theta.fix +ash</b>	0.07	0.09	0.08	0.13	0.13
<b>classtpx omega.fix</b>	0.09	0.13	0.11	0.22	0.19
<b>SVM (counts)</b>	0.09	0.22	0.20	0.21	0.15
<b>SVM (voom)</b>	0.05	0.05	0.03	0.08	0.05
<b>SVM (voom+ash)</b>	0.05	0.05	0.03	0.09	0.05
<b>PLS-LDA (counts)</b>	0.14	0.06	0.03	0.11	0.10
<b>PLS-LDA (voom)</b>	0.03	0.22	0.2	0.24	0.06
<b>PLS-LDA (voom +ash)</b>	0.03	0.16	0.16	0.18	0.06
<b>Normal</b>	0.09	0.12	0.09	0.18	0.11
<b>Poisson</b>	0.08	0.09	0.08	0.14	0.13
<b>Neg Binom</b>	0.10	0.18	0.10	0.23	0.17
<b>NN</b>	0.21	-	-	-	-

	All tissues	Arteries (3 tissues)	Arteries3000 (3 tissues)	Breast +Adipose	Breast + Adipose 3000
<b>classtpx theta.fix</b>	0.07	0.09	0.08	0.14	0.13
<b>classtpx theta.fix +ash</b>	0.07	0.09	0.08	0.13	0.13
<b>classtpx omega.fix</b>	0.09	0.13	0.11	0.22	0.19
<b>SVM (counts)</b>	0.09	0.22	0.20	0.21	0.15
<b>SVM (voom)</b>	0.05	0.05	0.03	0.08	0.05
<b>SVM (voom+ash)</b>	0.05	0.05	0.03	0.09	0.05
<b>PLS-LDA (counts)</b>	0.14	0.06	0.03	0.11	0.10
<b>PLS-LDA (voom)</b>	0.03	0.22	0.2	0.24	0.06
<b>PLS-LDA (voom +ash)</b>	0.03	0.16	0.16	0.18	0.06
<b>Normal</b>	0.09	0.12	0.09	0.18	0.11
<b>Poisson</b>	0.08	0.09	0.08	0.14	0.13
<b>Neg Binom</b>	0.10	0.18	0.10	0.23	0.17
<b>NN</b>	0.21	-	-	-	-

How about single-cell data?

**STRUCTURE by development phase**

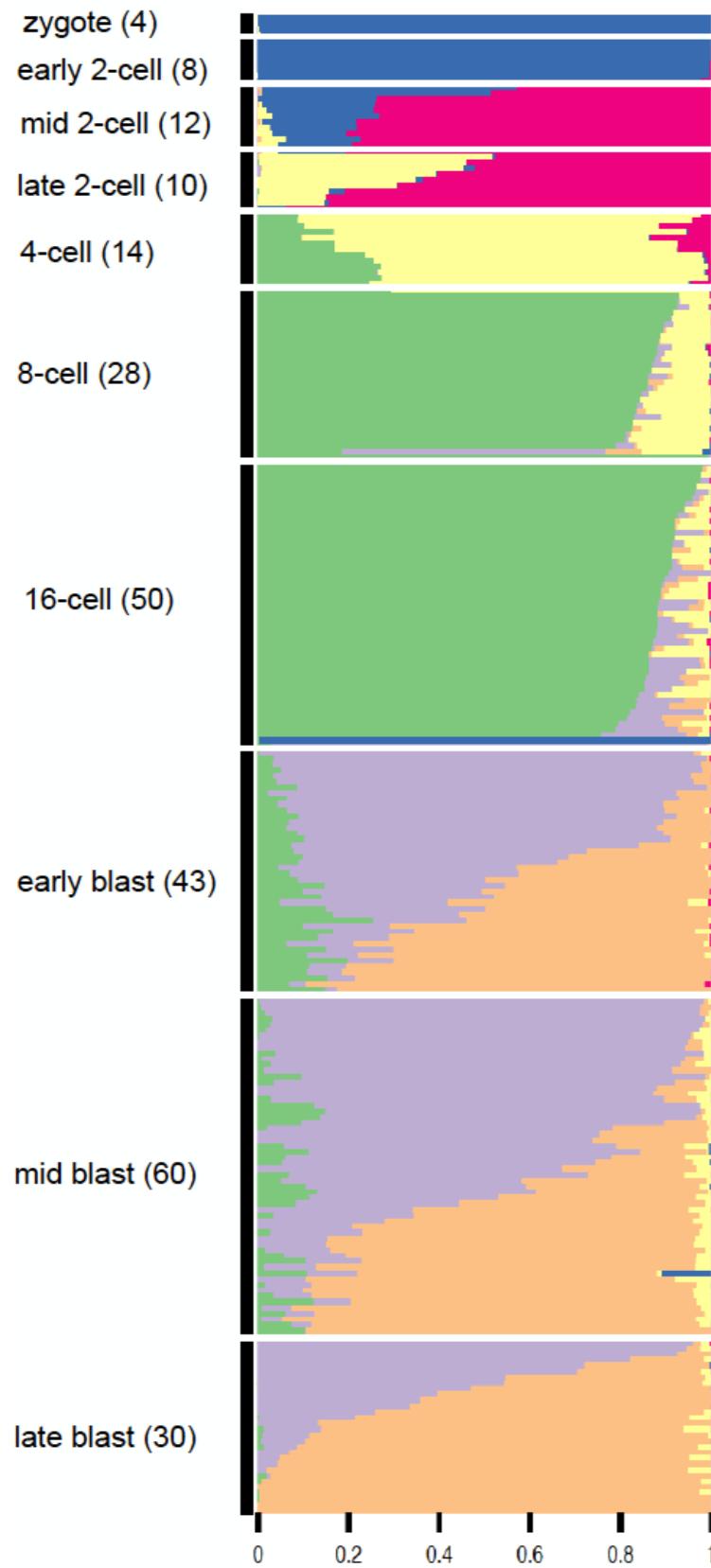


*Deng et al (2014)  
single cell  
development phase  
data*

maptpx  
clustering

6 clusters

## STRUCTURE by development phase

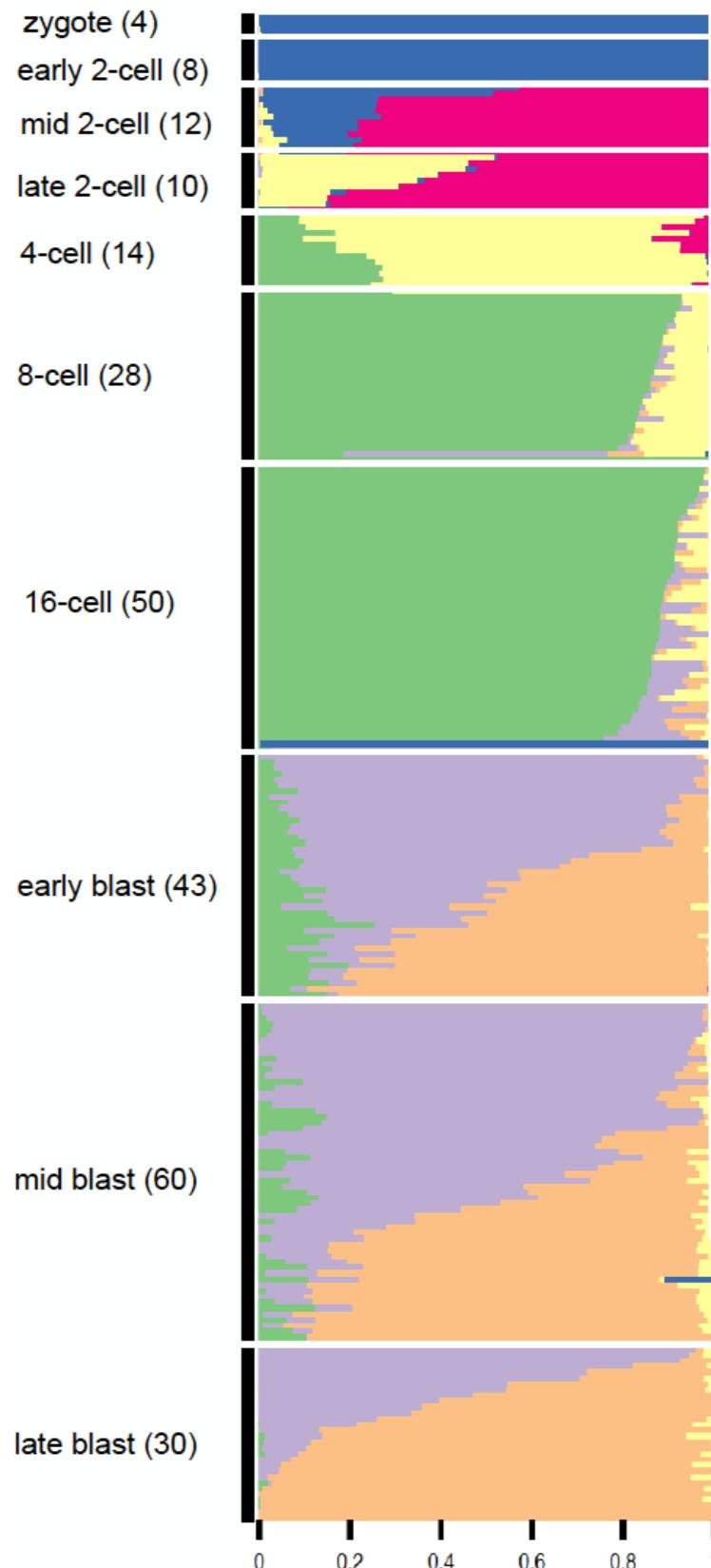


**Deng et al (2014)  
single cell  
development phase  
data**

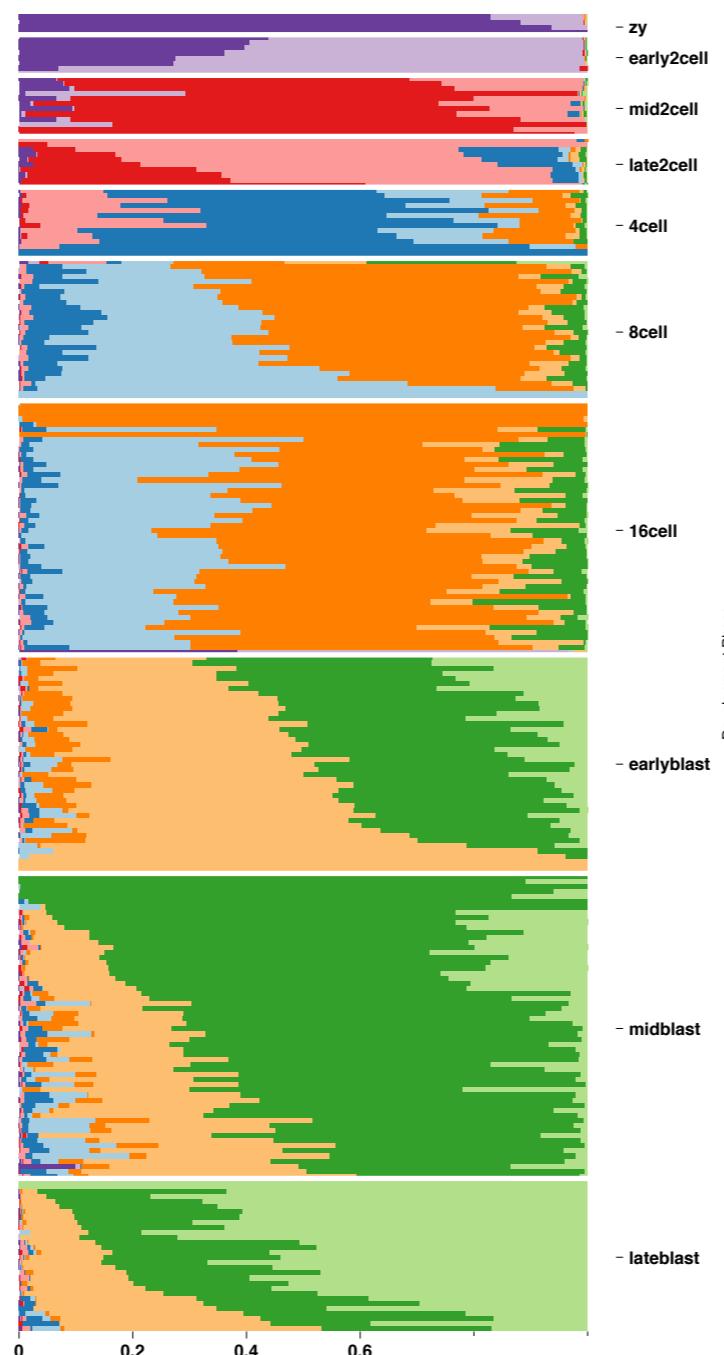
We pick 20 % samples from  
each phase as training  
for classtpx and other classifiers

6 clusters

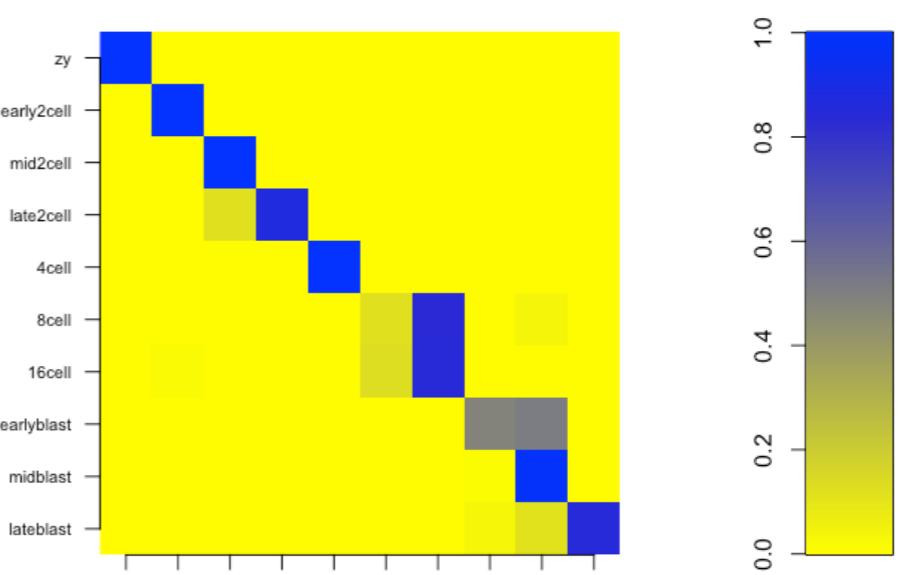
## STRUCTURE by development phases



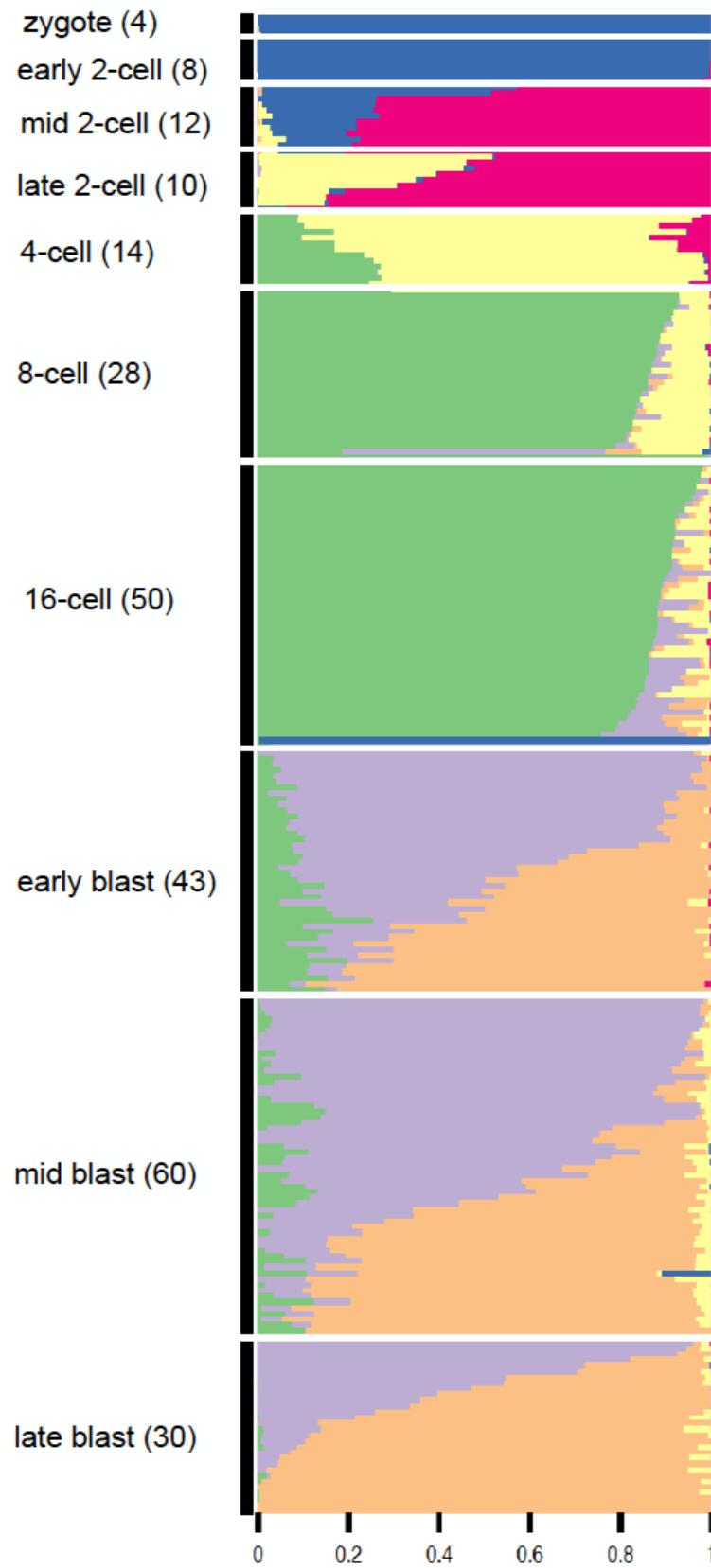
6 clusters



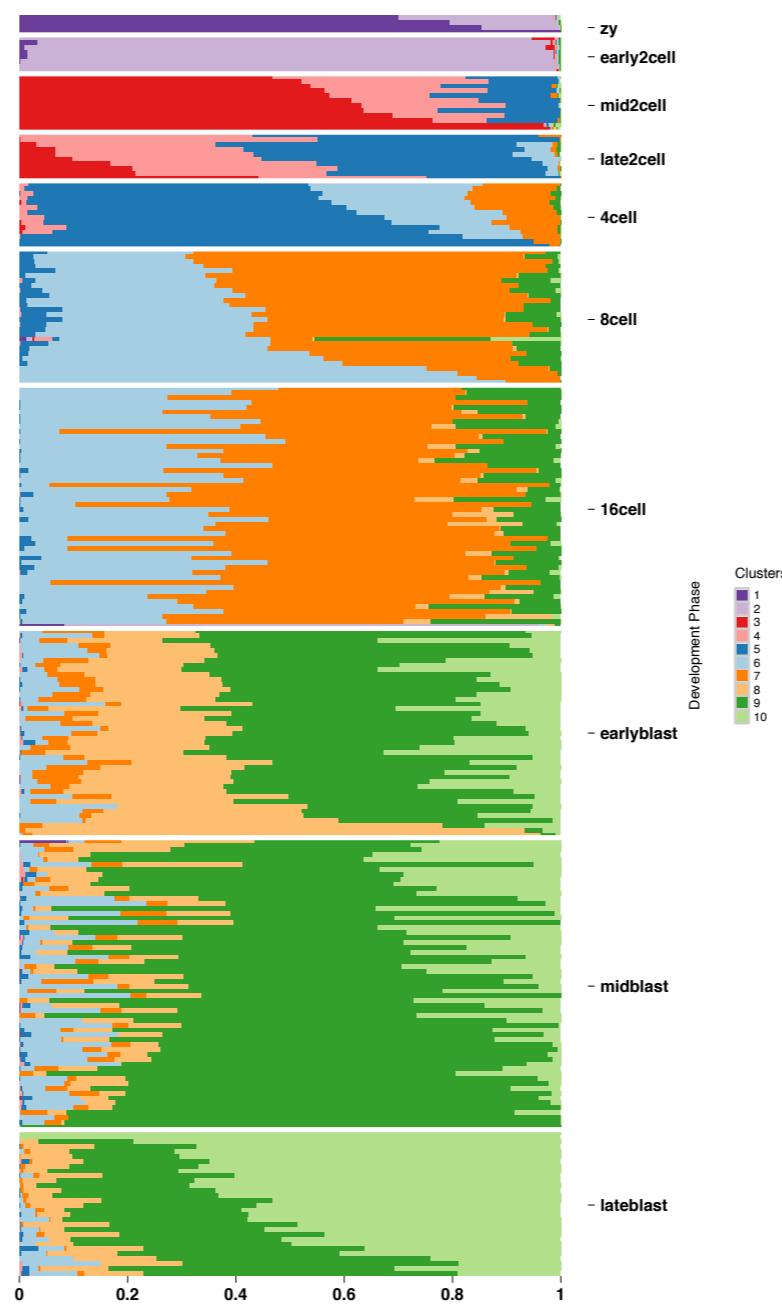
**theta.fix**



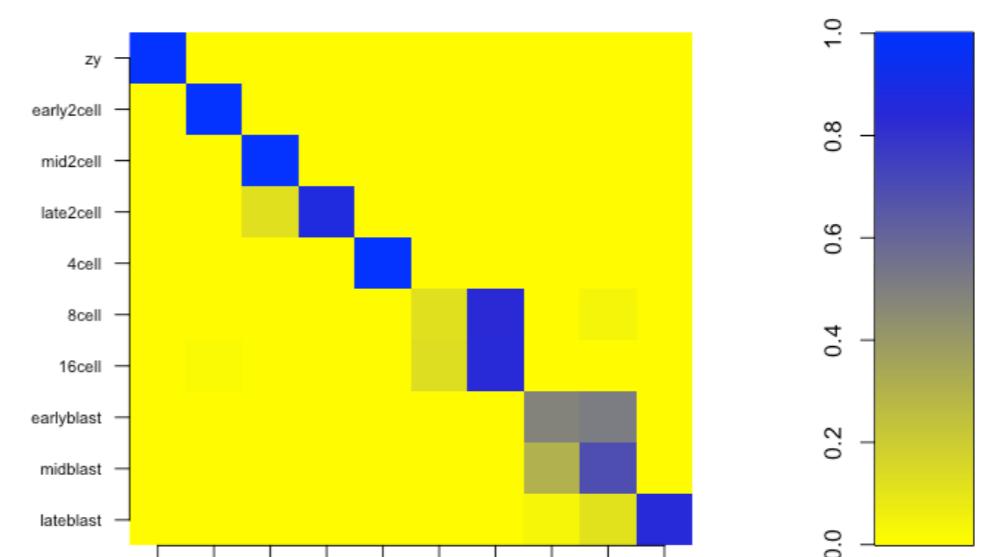
## STRUCTURE by development phase



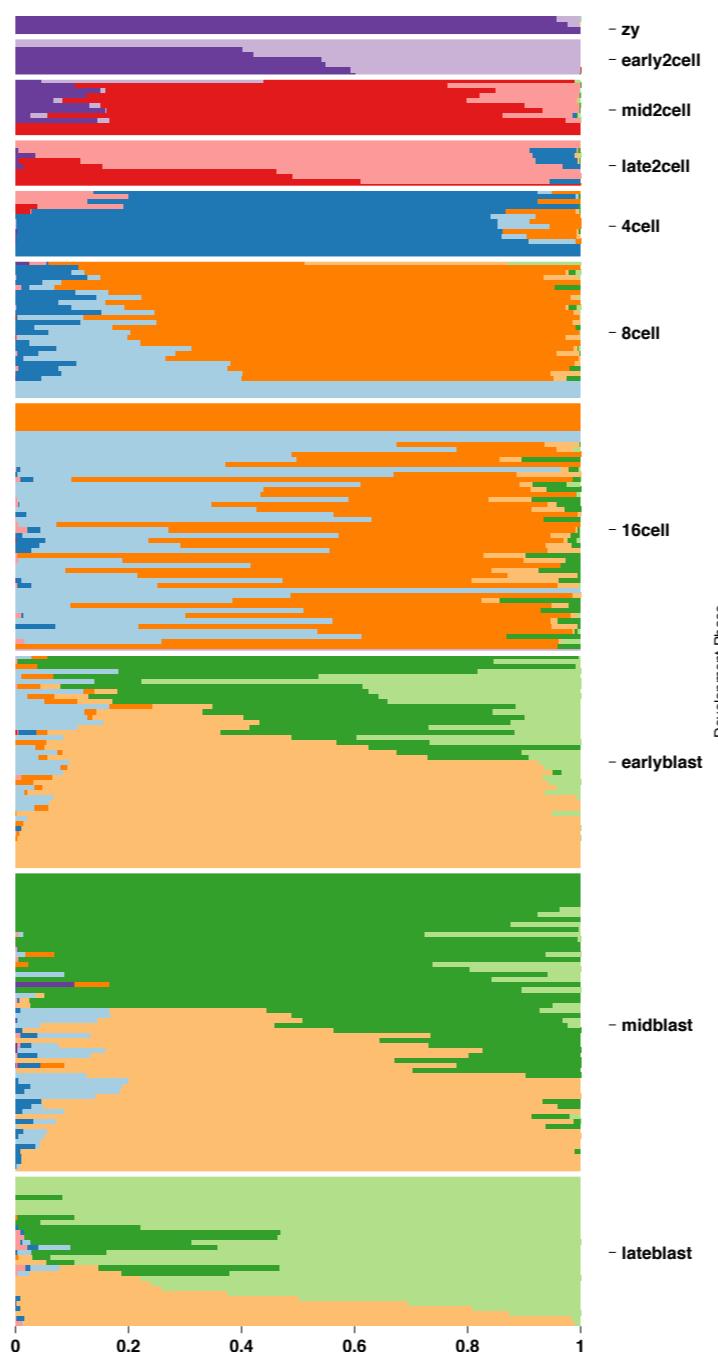
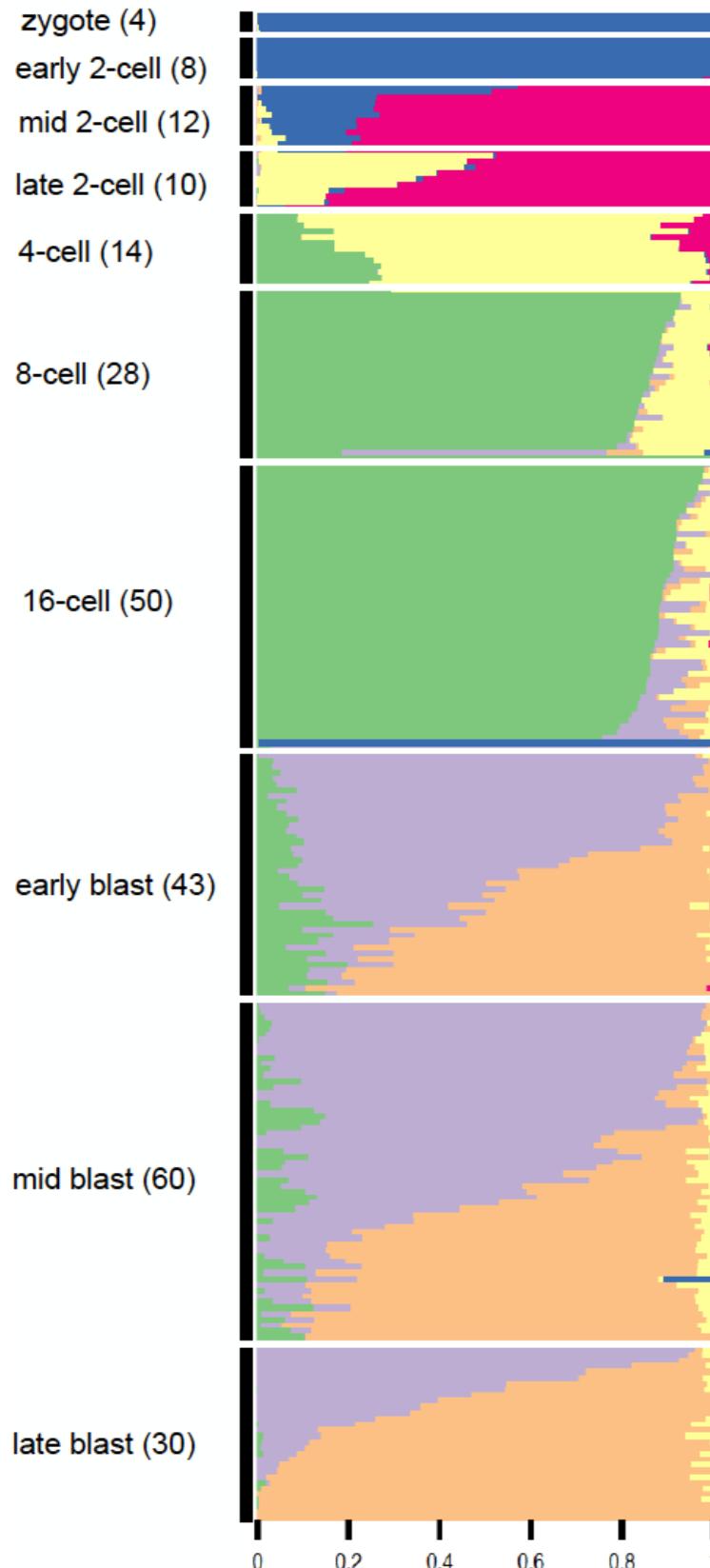
6 clusters



**theta.fix+ash**



## STRUCTURE by development phas



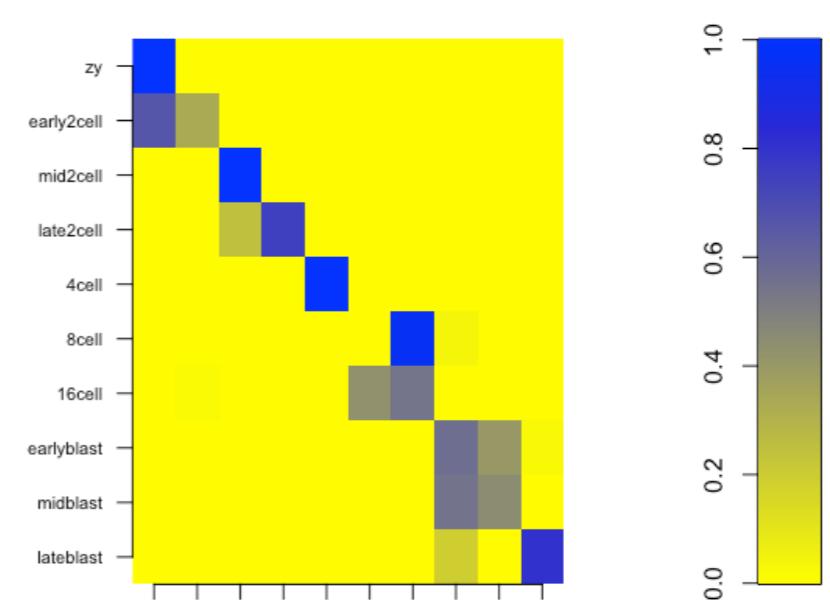
Clusters

- 1
- 2
- 4
- 5
- 6
- 7
- 8
- 9
- 10

Development Phase

**omega.fix**

6 clusters



	Deng et al (2014)
classtpx theta.fix	0.24
classtpx theta.fix +ash	0.31
classtpx omega.fix	0.40
SVM (counts)	> 0.6
SVM (voom)	0.56
SVM (voom+ash)	> 0.6
PLS-LDA (counts)	0.43
PLS-LDA (voom)	0.38
PLS-LDA (voom+ash)	0.59
Normal	> 0.6
Poisson	0.32
Neg Binom	0.30

# **Summary**

classtpx seems to be robust to the size of the training sample, unlike other classifiers (SVM, PLS-LDA etc) as in Deng et al (2014)

theta.fix and theta.fix+ash seem to perform better as classifier than omega.fix

As a classifier, SVM on voom or voom+ash data seems to outperform classtpx, but not by a big margin

classtpx performs soft clustering- assigns grade of membership of test data to each class, other methods not meant to serve that purpose.