

Dimensionality reduction by UMAP  
reinforces sample heterogeneity analysis in  
bulk transcriptomic data - Supplementary

Table S2: Dimensionality reduction methods, 4 methods

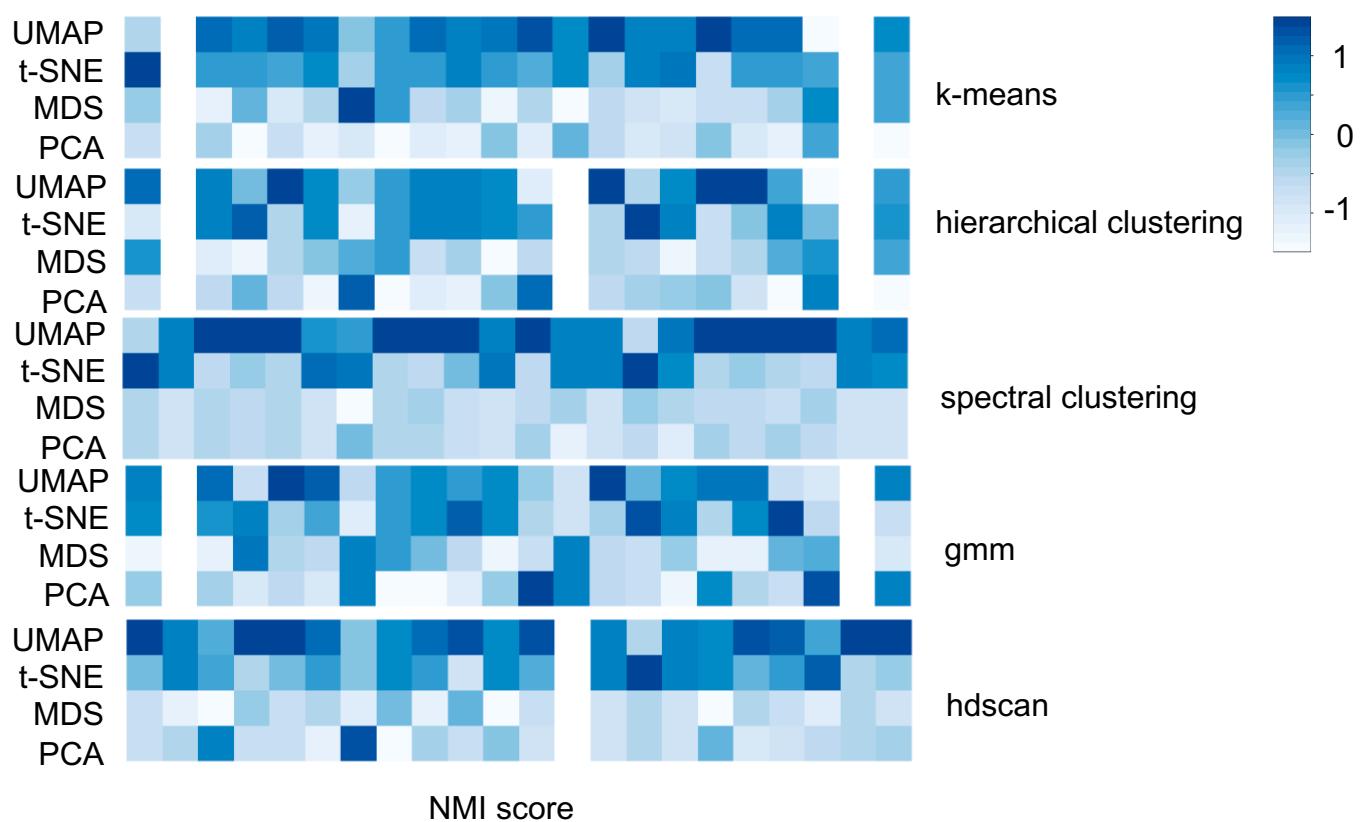
Algorithm	Version	Availability	Implementation language	Parameter
UMAP	0.3.10	<a href="https://github.com/lmc-innes/umap/">https://github.com/lmc-innes/umap/</a>	Python	<code>n_components=2, n_neighbors=15, min_dist=0.1, metric='euclidean'</code>
t-SNE	scikit-learn 0.23.1	<a href="https://scikit-learn.org/">https://scikit-learn.org/</a>	Python	<code>n_components=2, perplexity=30.0, metric='euclidean', early_exaggeration=12.0, learning_rate=200.0, n_iter=1000,</code>
MDS	scikit-learn 0.23.1	<a href="https://scikit-learn.org/">https://scikit-learn.org/</a>	Python	<code>n_components=2, *, metric=True, n_init=4, max_iter=300, verbose=0, eps=0.001, n_jobs=None, random_state=None, dissimilarity='euclidean'</code>
PCA	scikit-learn 0.23.1	<a href="https://scikit-learn.org/">https://scikit-learn.org/</a>	Python	<code>n_components=2, copy=True, whiten=False, svd_solver='auto', tol=0.0, iterated_power='auto', random_state=None</code>

Table S3: 5 clustering algorithms

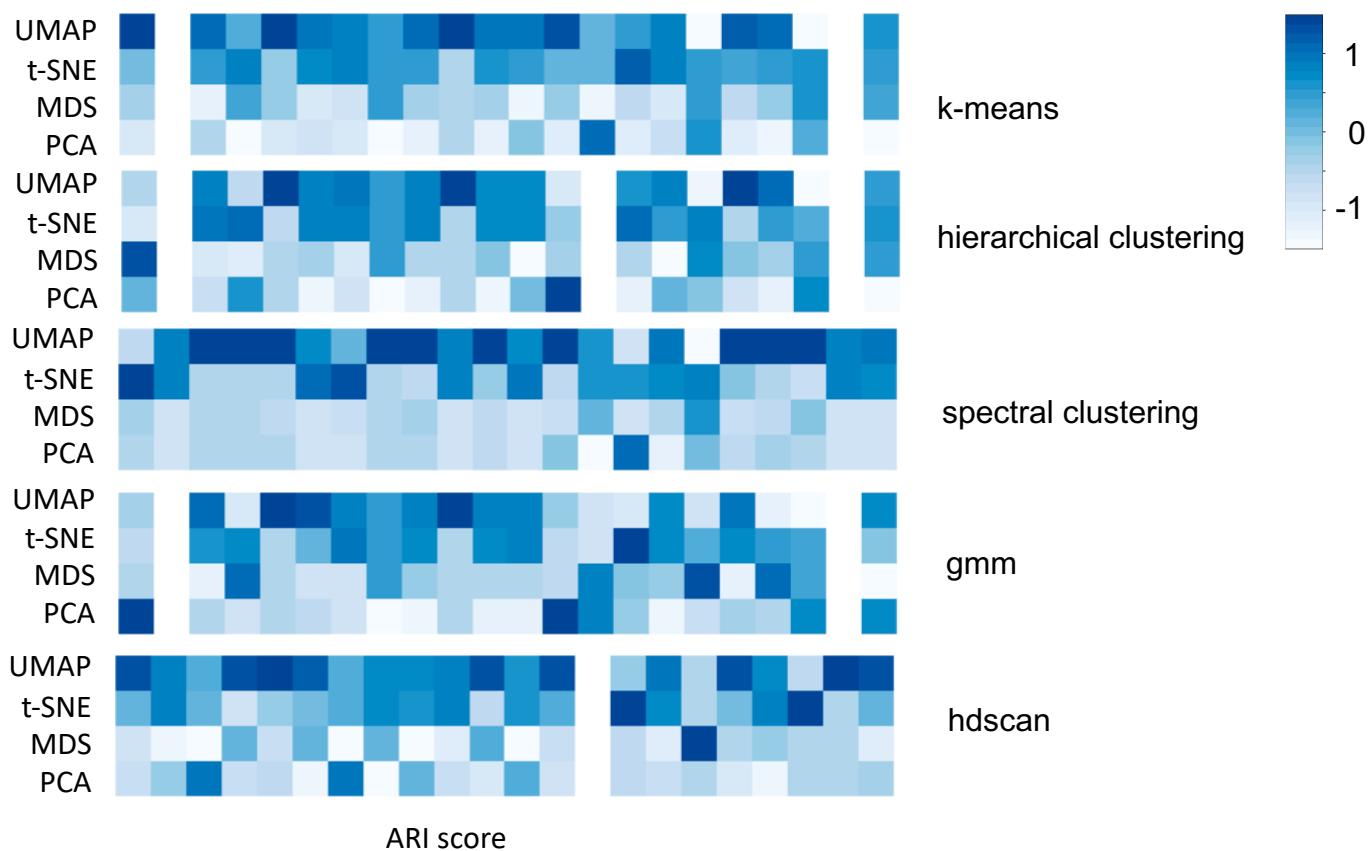
Algorithm	Version	Availability	Implementation language
k-means	scikit-learn 0.23.1: sklearn.cluster.KMeans	<a href="https://scikit-learn.org/">https://scikit-learn.org/</a>	Python
Hierarchical clustering	scikit-learn 0.23.1: sklearn.cluster.AgglomerativeClustering	<a href="https://scikit-learn.org/">https://scikit-learn.org/</a>	Python
Spectral clustering	scikit-learn 0.23.1: sklearn.cluster.SpectralClustering	<a href="https://scikit-learn.org/">https://scikit-learn.org/</a>	Python
gmm	scikit-learn 0.23.1: sklearn.mixture.GaussianMixture	<a href="https://scikit-learn.org/">https://scikit-learn.org/</a>	Python
hdbscan	hdbscan 0.8.18	<a href="https://hdbscan.readthedocs.io/">https://hdbscan.readthedocs.io/</a>	Python

Figure S1: Clustering accuracy of five clustering algorithms on embedded space by four dimensionality reduction methods

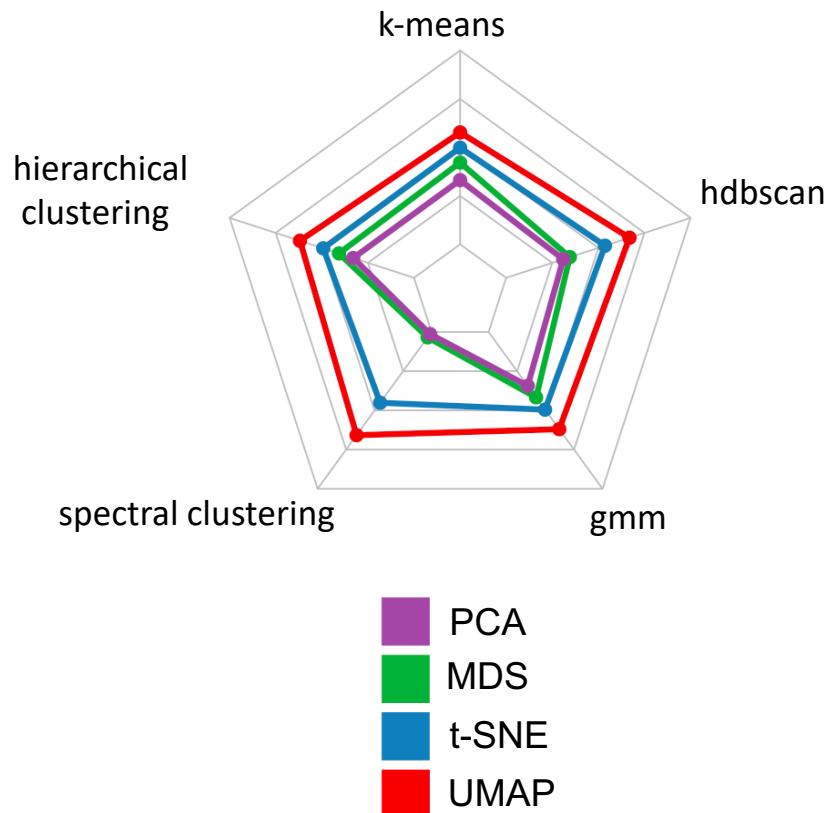
a



b

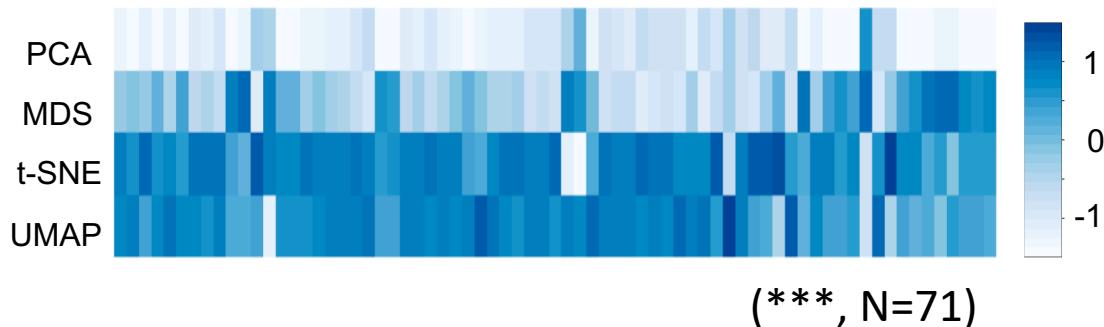


C



- Heatmap for clustering accuracy by normalized mutual information (NMI). Darker blue means higher clustering accuracy. Five clustering algorithms are listed on the right side. [“1” in the heatmap represent the perfect accuracy, “0”.... “-1” represent....]
- Heatmap for clustering accuracy by adjusted Rand index (ARI).
- Radar plot of clustering accuracy (ARI score) comparison using five clustering methods. The average ARI score was on 22 datasets with cluster labels. The input was the embedded two-dimensional coordinates of each dimensionality reduction methods. Larger scale denotes better clustering accuracy, and UMAP outperformed the other three.

Figure S2: Neighborhood preserving (knn\_k = 30), Average Jaccard index with 30 neighborhoods.



Heatmap for evaluating neighbourhood preserving of each method on 71 datasets. The number of neighbours is set as 30. The darker the colour is, the better the local information is retained.

Figure S3: Visualization of dataset GSE98793 and GSE107990 showing batch effects in two-dimensional space by dimensionality reduction methods.

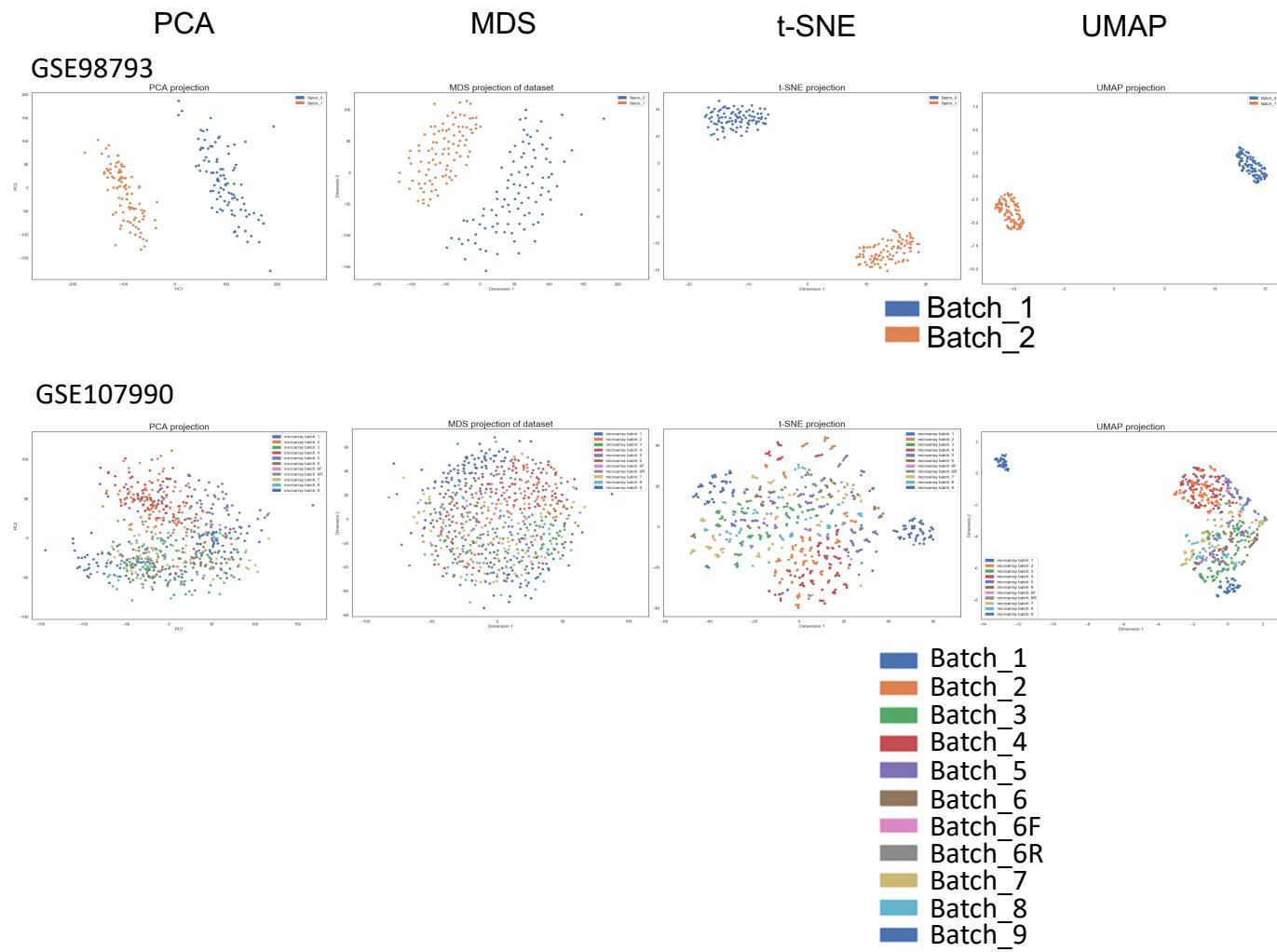
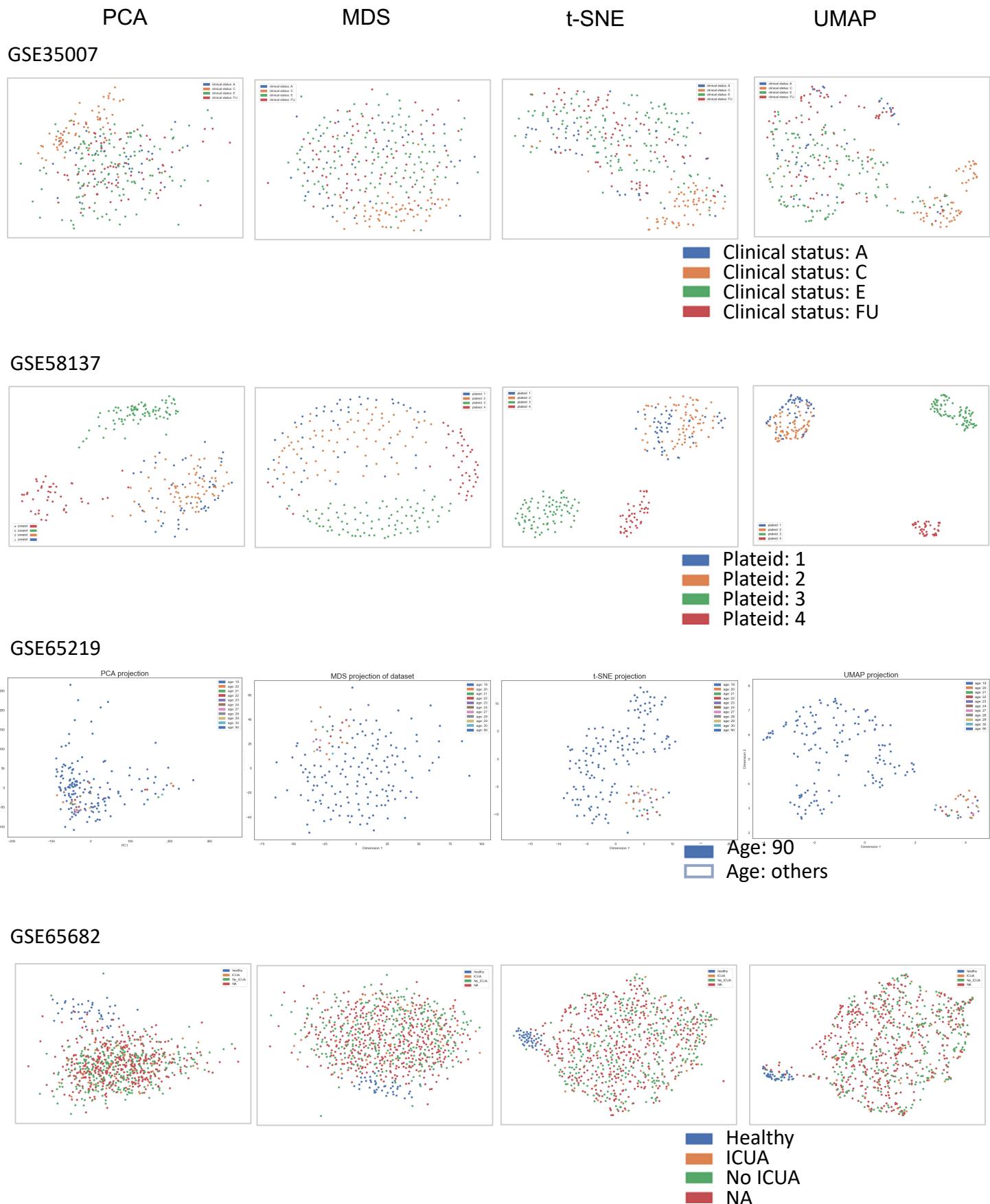


Figure S4: Visualization of eight datasets illustrating biological groups by dimensionality reduction methods. Colours represents groups.



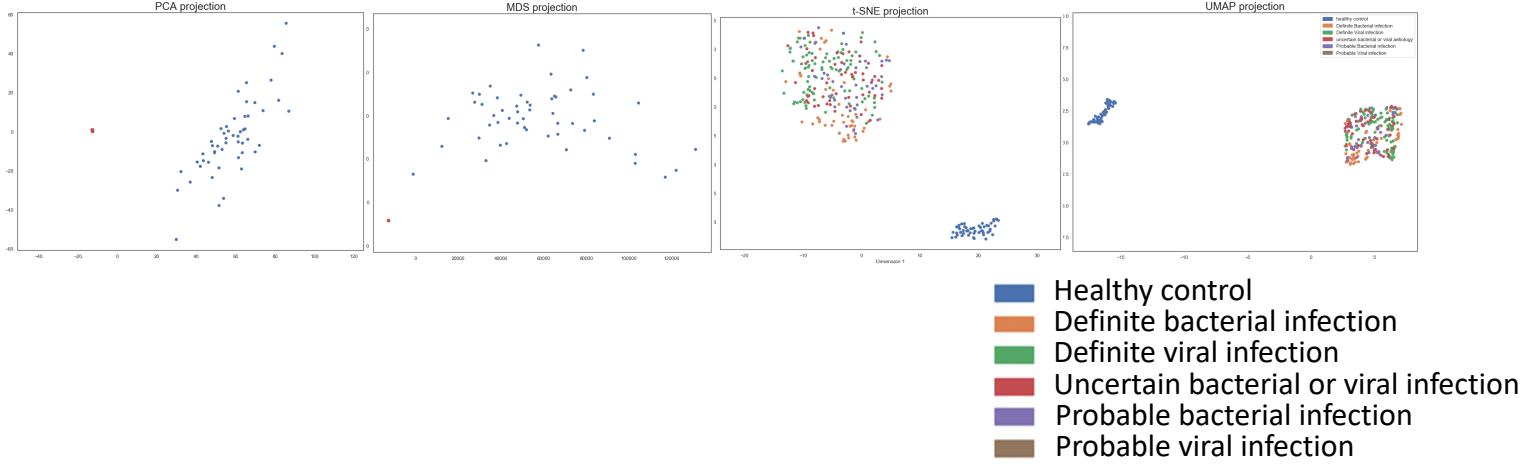
## PCA

## MDS

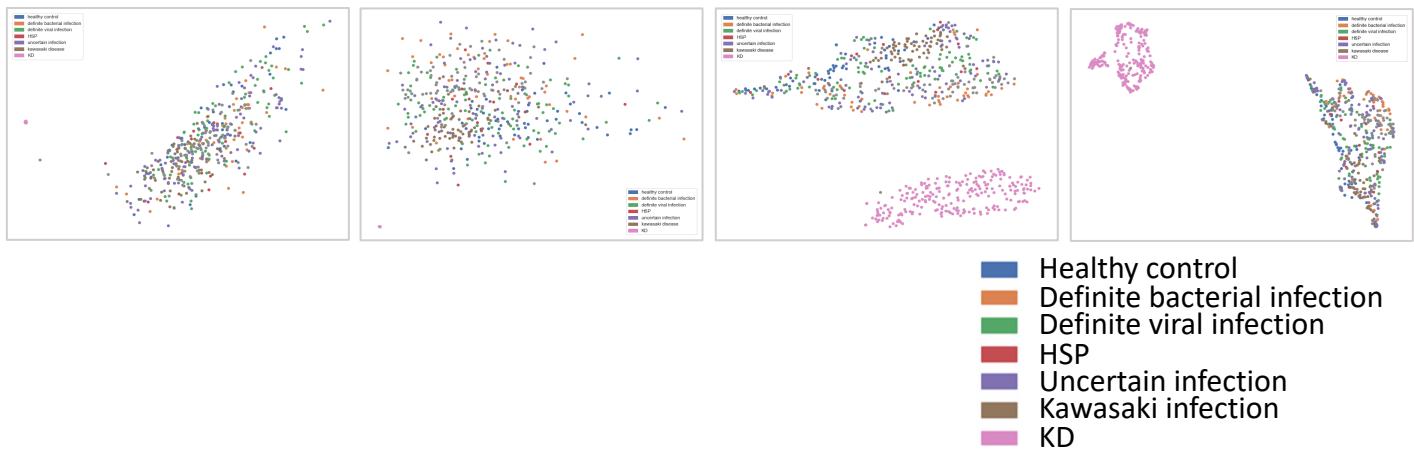
## t-SNE

## UMAP

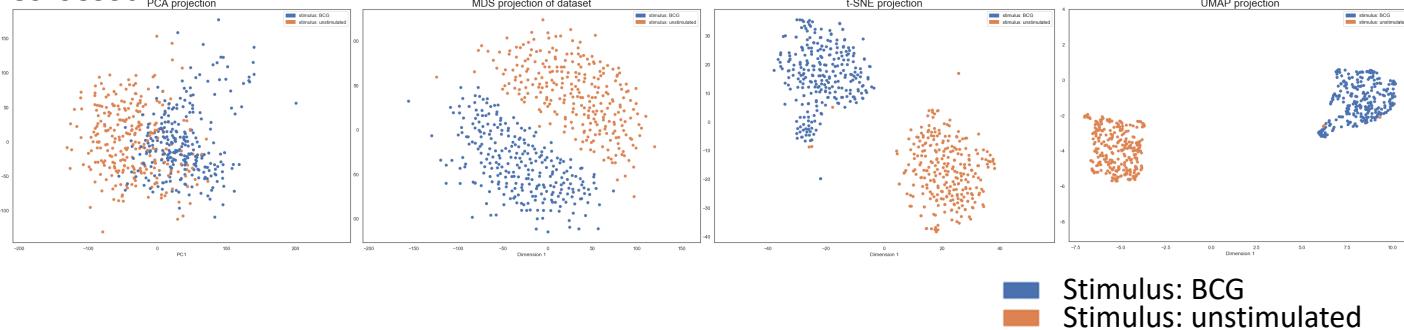
GSE72809



GSE73464



GSE98550



GSE111368

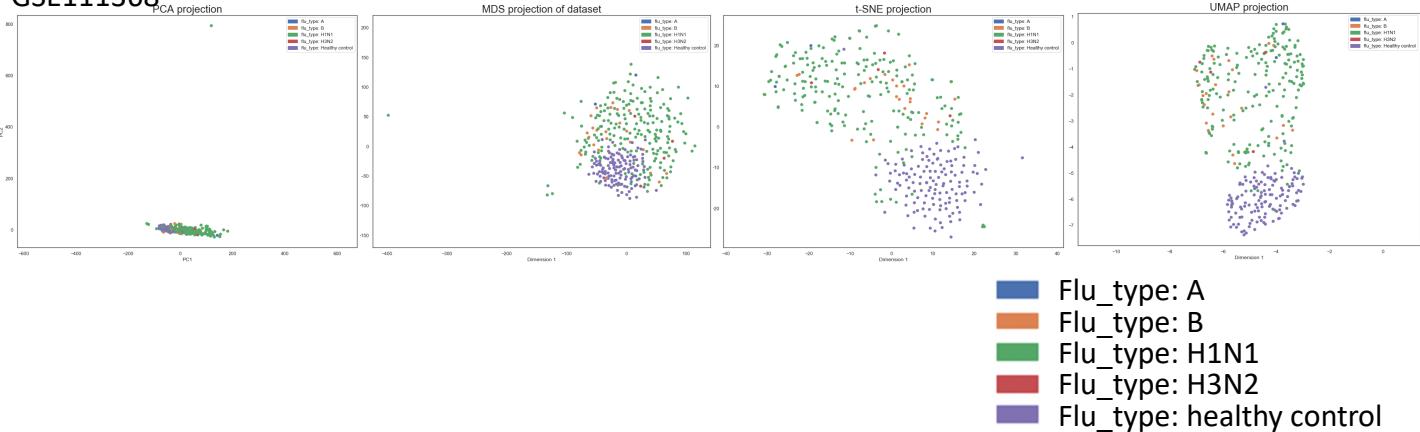
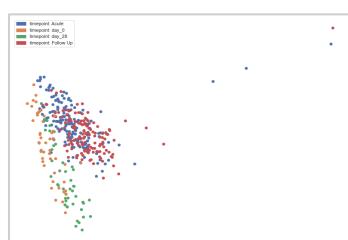


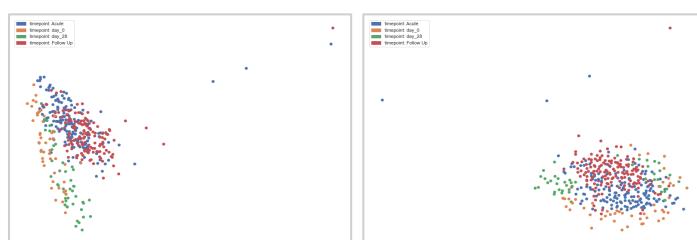
Figure S5: Associating sample features to clustering structure

PCA

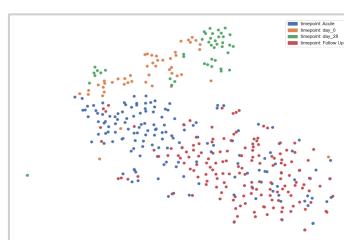
GSE61821



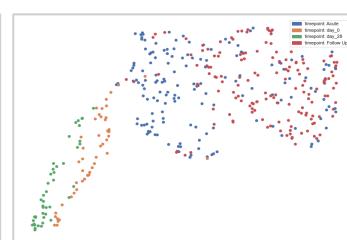
MDS



t-SNE

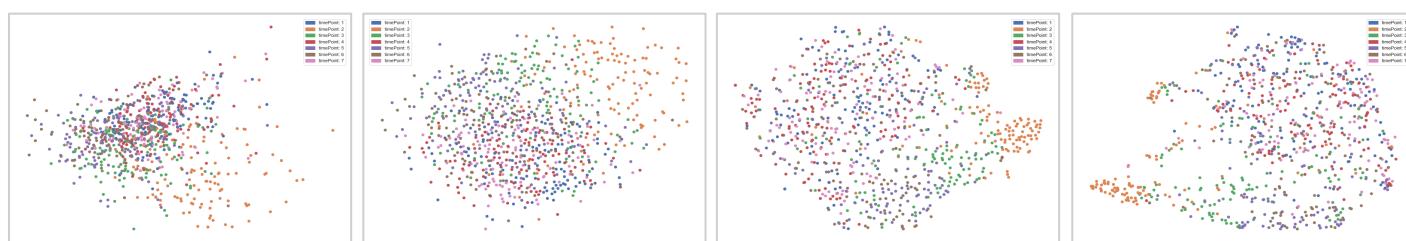


UMAP



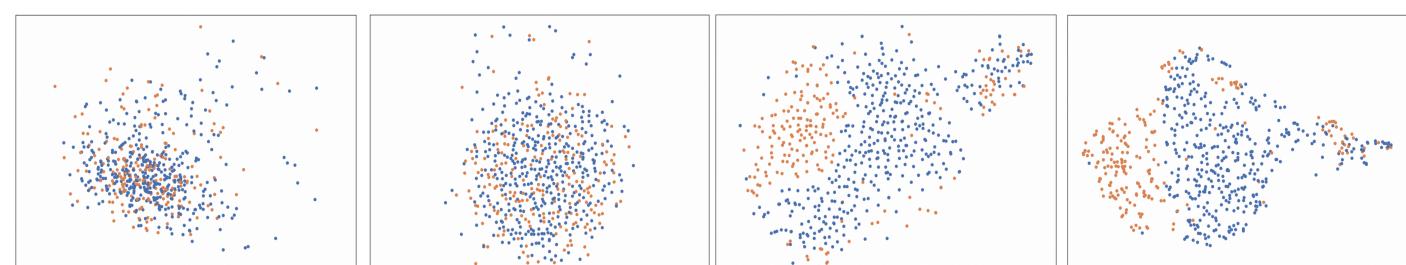
Timepoint: acute  
Timepoint: day\_0  
Timepoint: day\_28  
Timepoint: follow up

GSE68310



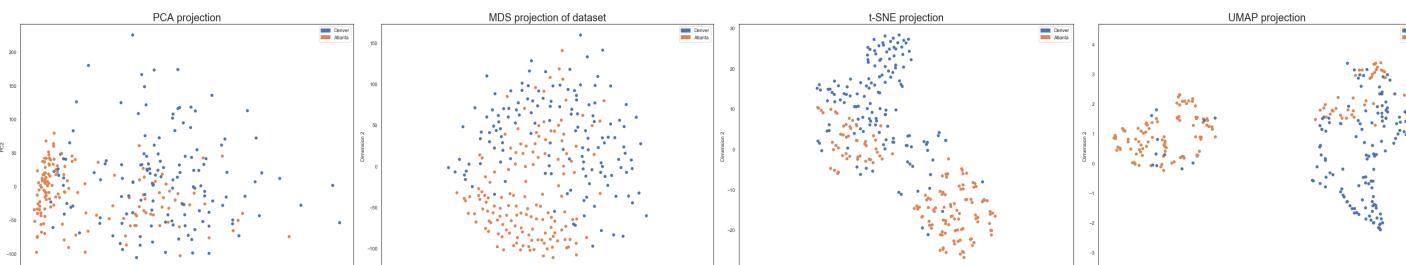
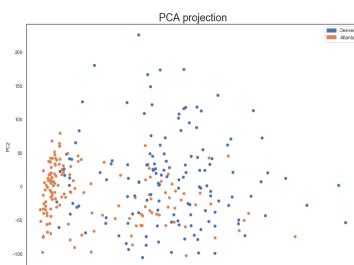
timepoint: 1  
timepoint: 2  
timepoint: 3  
timepoint: 4  
timepoint: 5  
timepoint: 6  
timepoint: 7

GSE71220



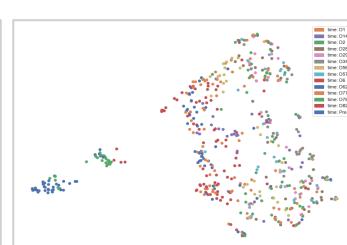
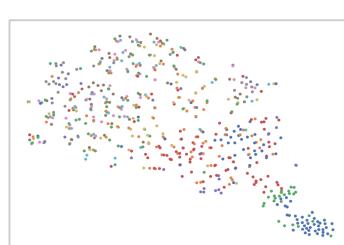
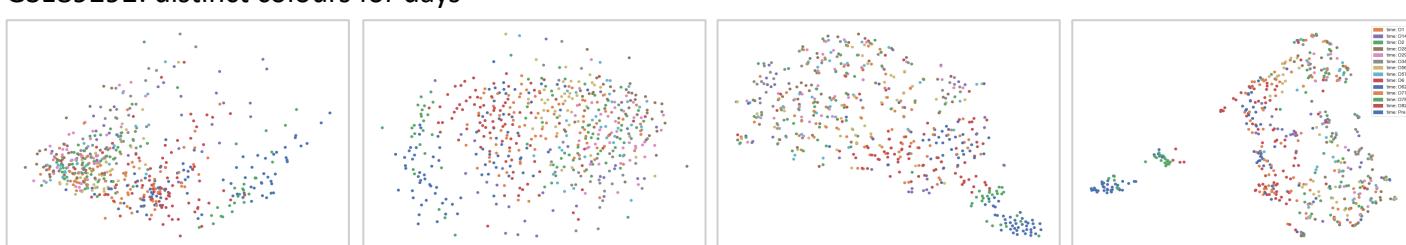
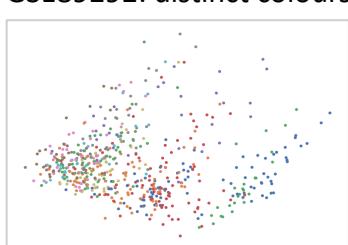
Male  
Female

GSE79396



Place: Denver  
Place: Atlanta

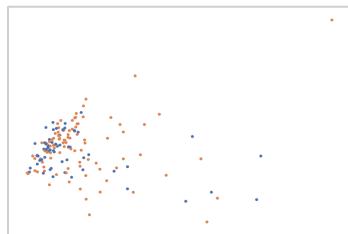
GSE89292: distinct colours for days



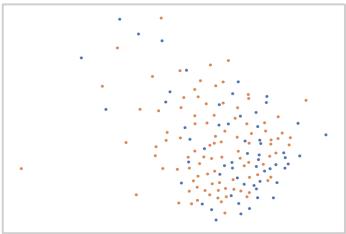
Day 014  
Day 015  
Day 016  
Day 017  
Day 018  
Day 019  
Day 020  
Day 021  
Day 022  
Day 023

PCA

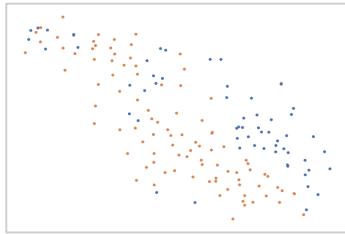
GSE110551



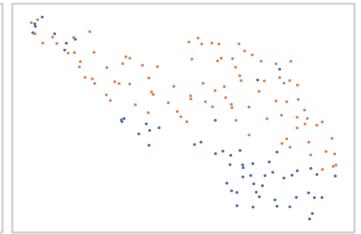
MDS



t-SNE

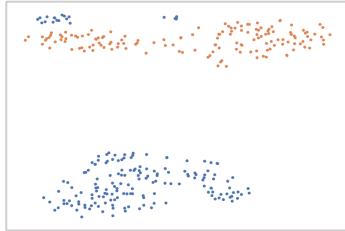
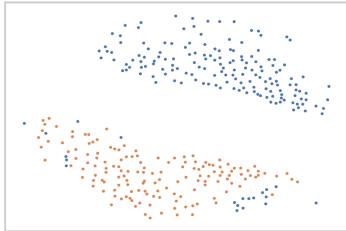


UMAP



Male  
Female

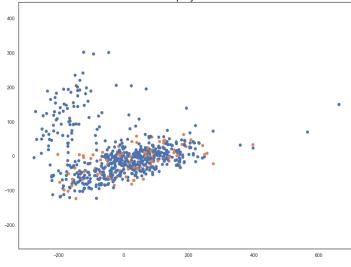
GSE113867



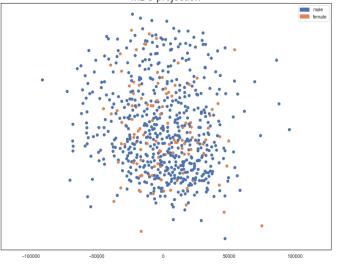
Date: April-11-2018  
Date: April-30-2018

GSE125216

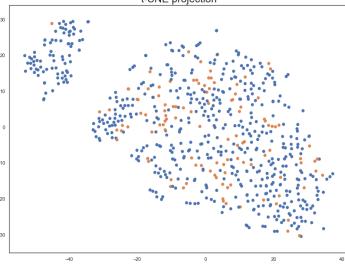
PCA projection



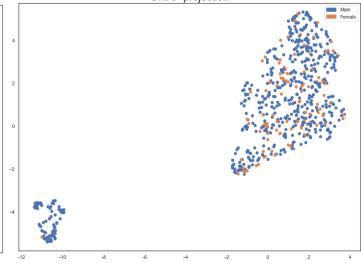
MDS projection



t-SNE projection

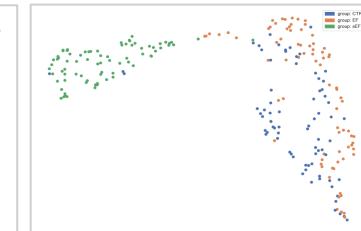
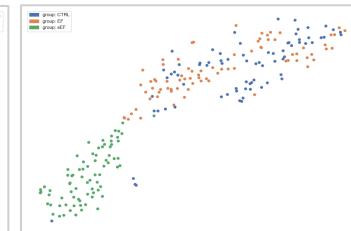
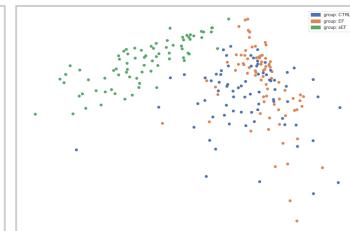
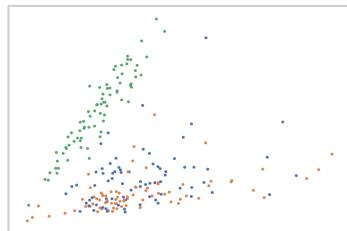


UMAP projection



Male  
Female

GSE133822



CTRL  
EF  
SEF

Figure S6: Visualization of datasets showing clustering structures associated no pre-defined features

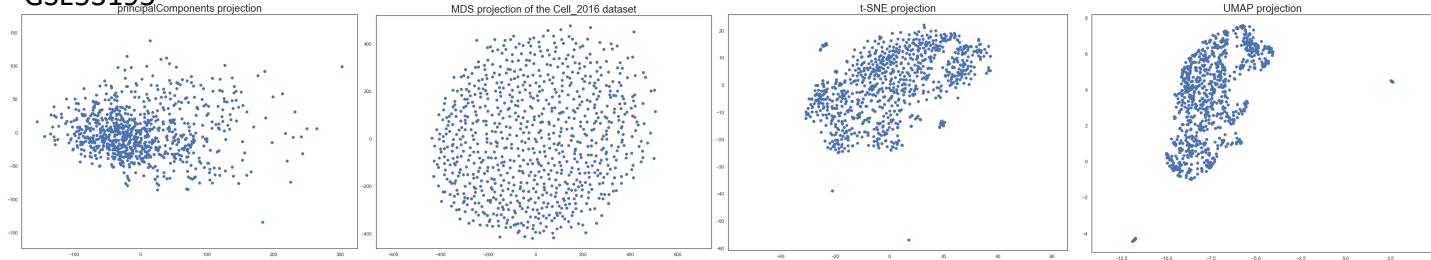
PCA

MDS

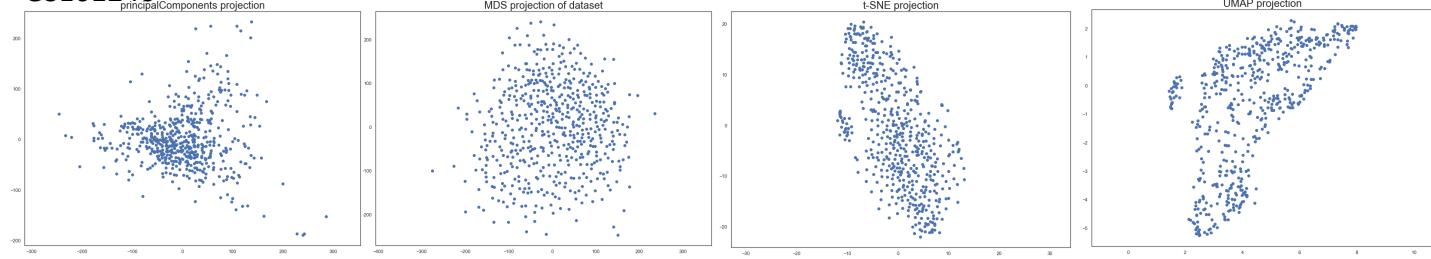
t-SNE

UMAP

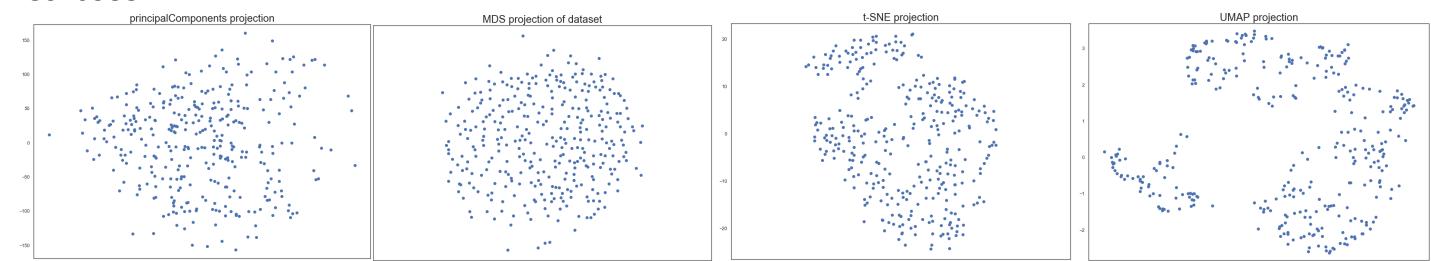
GSE53195



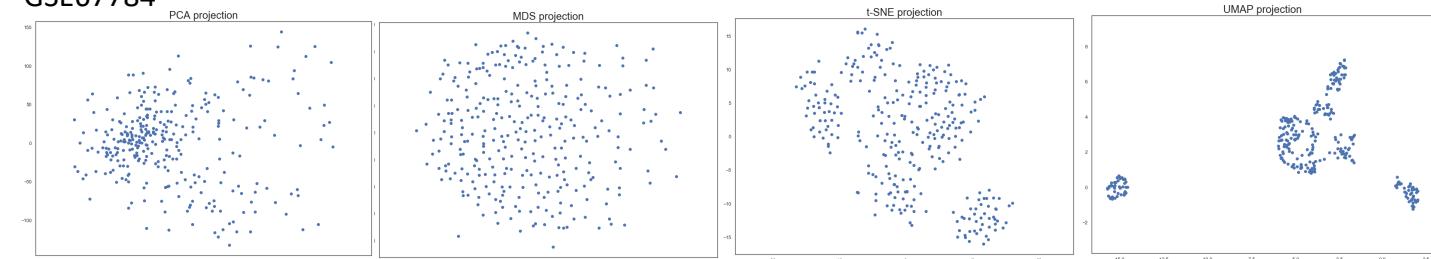
GSE61240



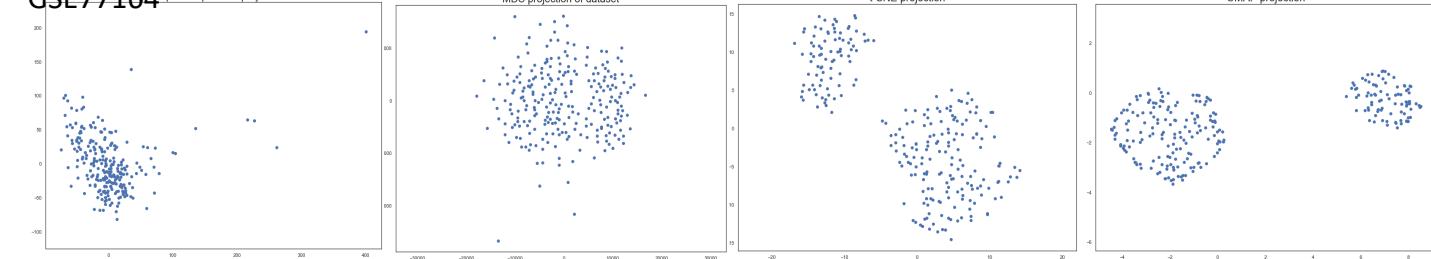
GSE63881



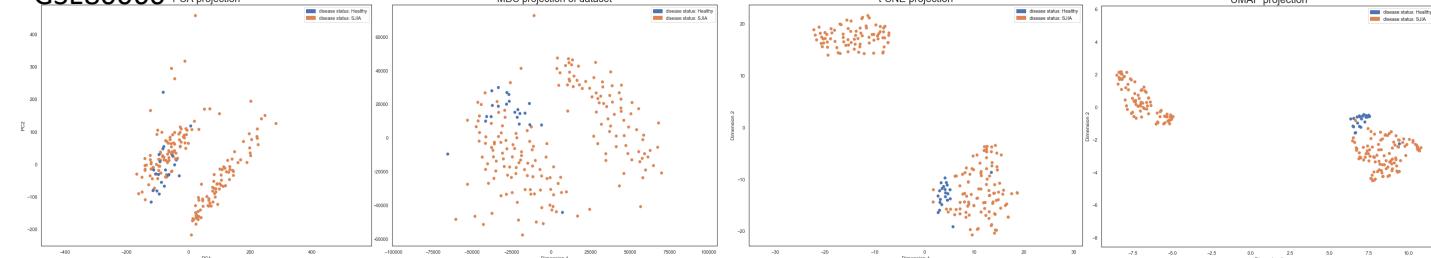
GSE67784



GSE77164



GSE80060



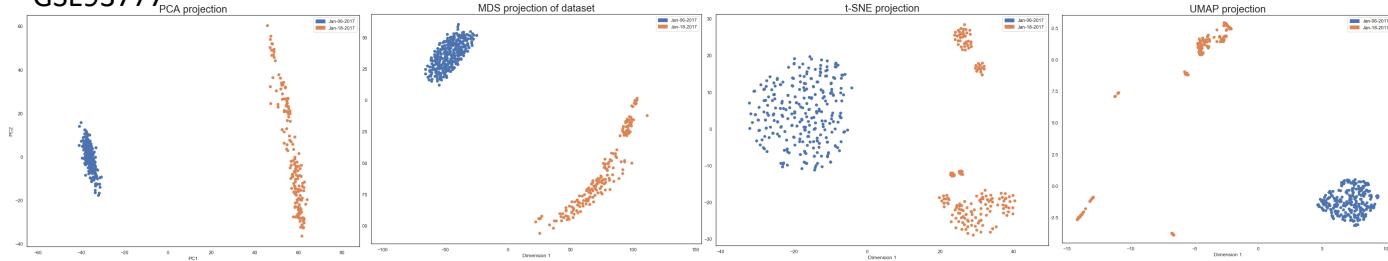
PCA

GSE93777

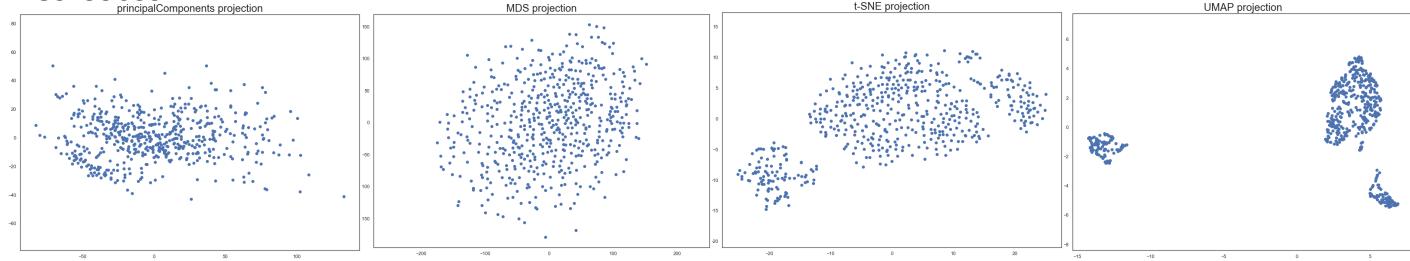
MDS

t-SNE

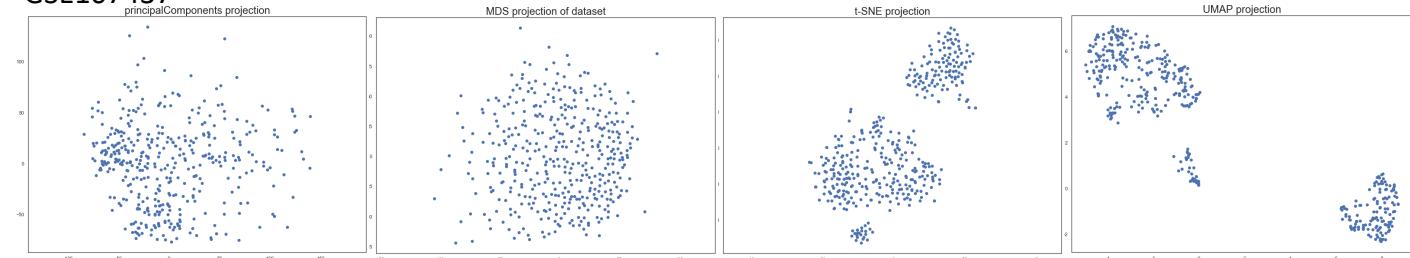
UMAP



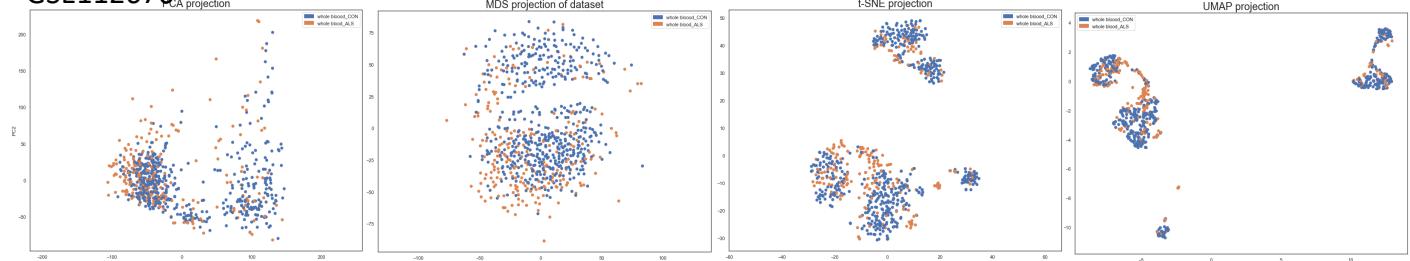
GSE99039



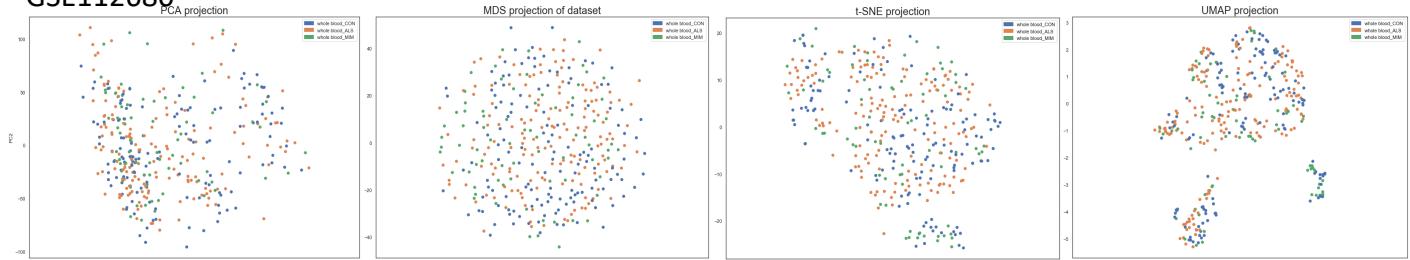
GSE107437



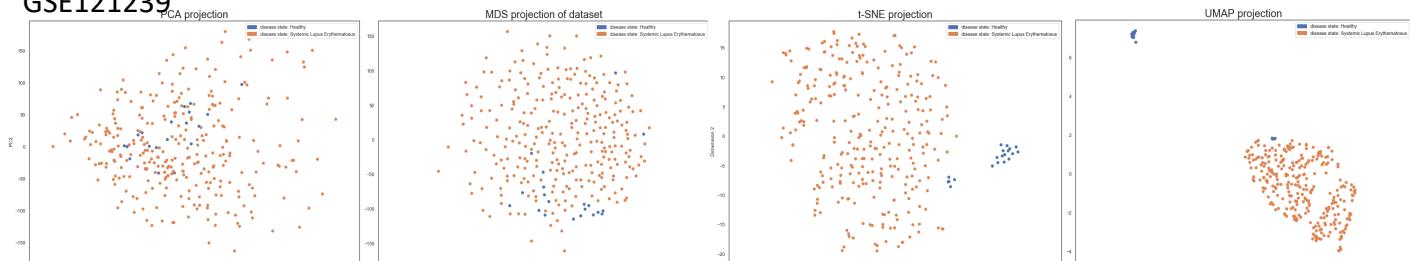
GSE112676



GSE112680

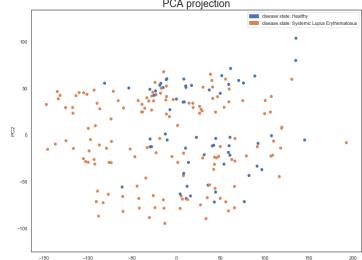


GSE121239



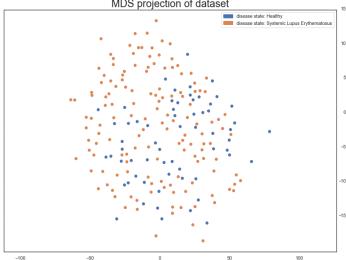
# PCA

GSE130953



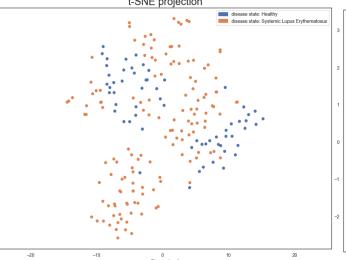
# MDS

MDS projection of dataset



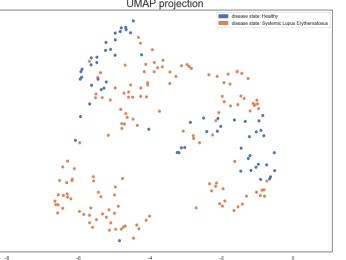
# t-SNE

t-SNE projection

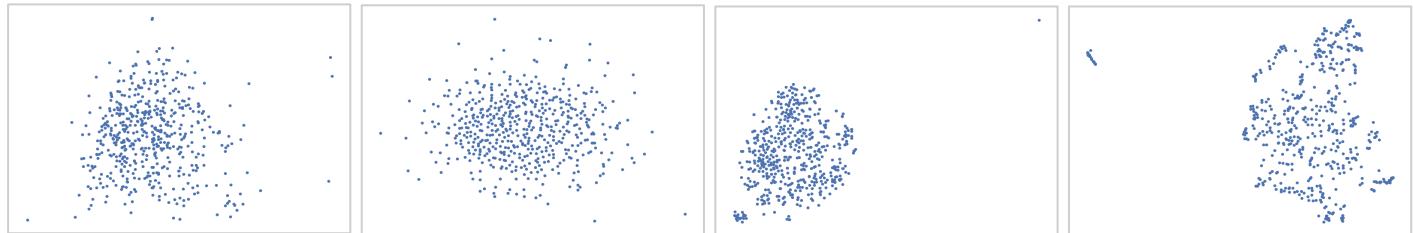


# UMAP

UMAP projection



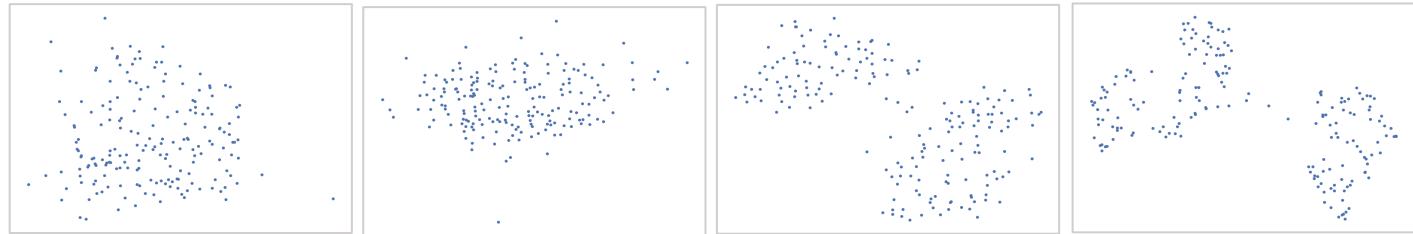
GSE47755



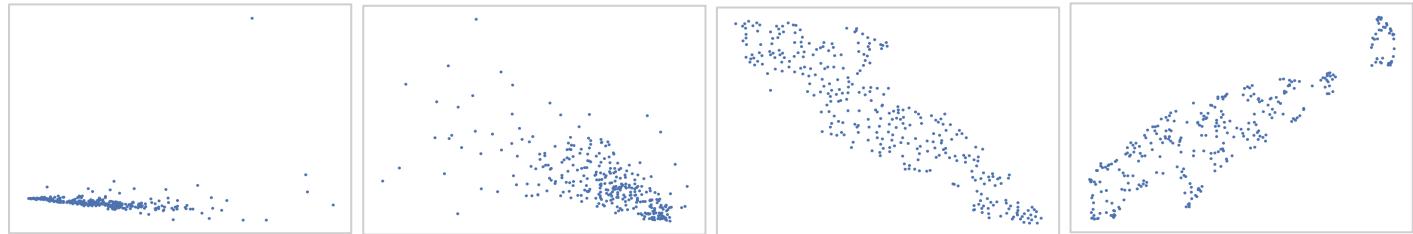
GSE64930



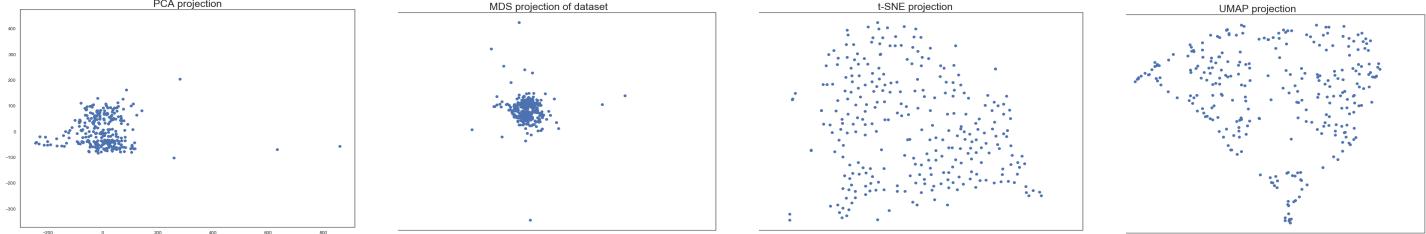
GSE85531



GSE97356



GSE124284



GSE124326



GSE124400

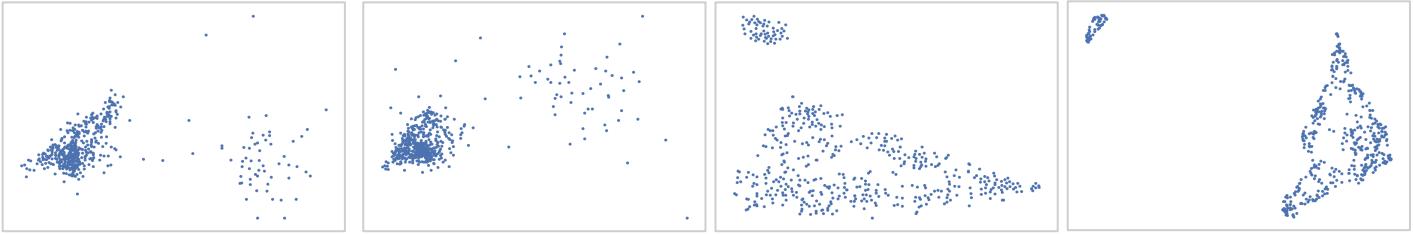


Figure S7: Gene set enrichment analysis between sG1 v.s. sG0 and sG2 v.s. sG0 with top 20 differentially regulated molecular pathways ranked by adjusted p-value.

