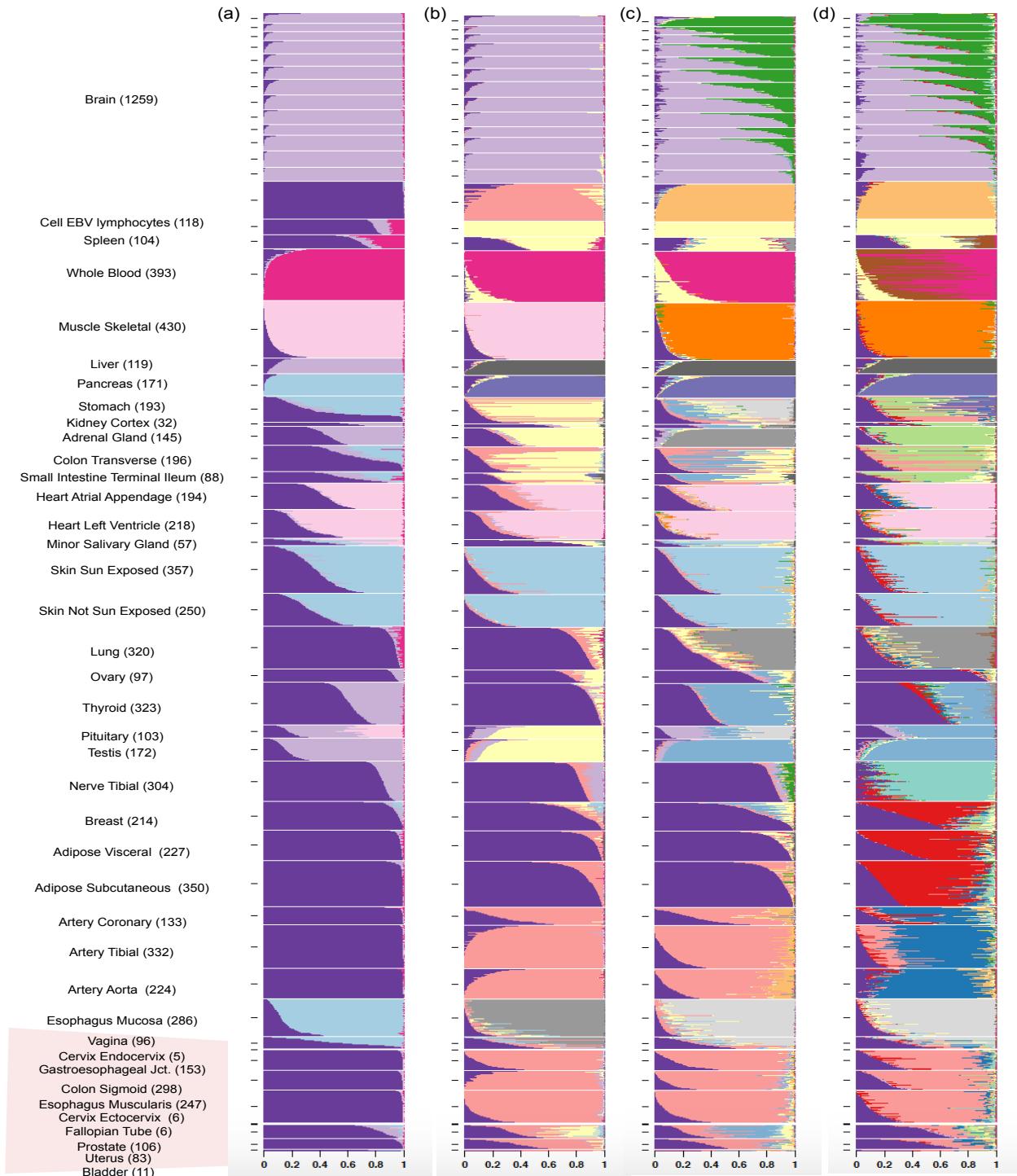
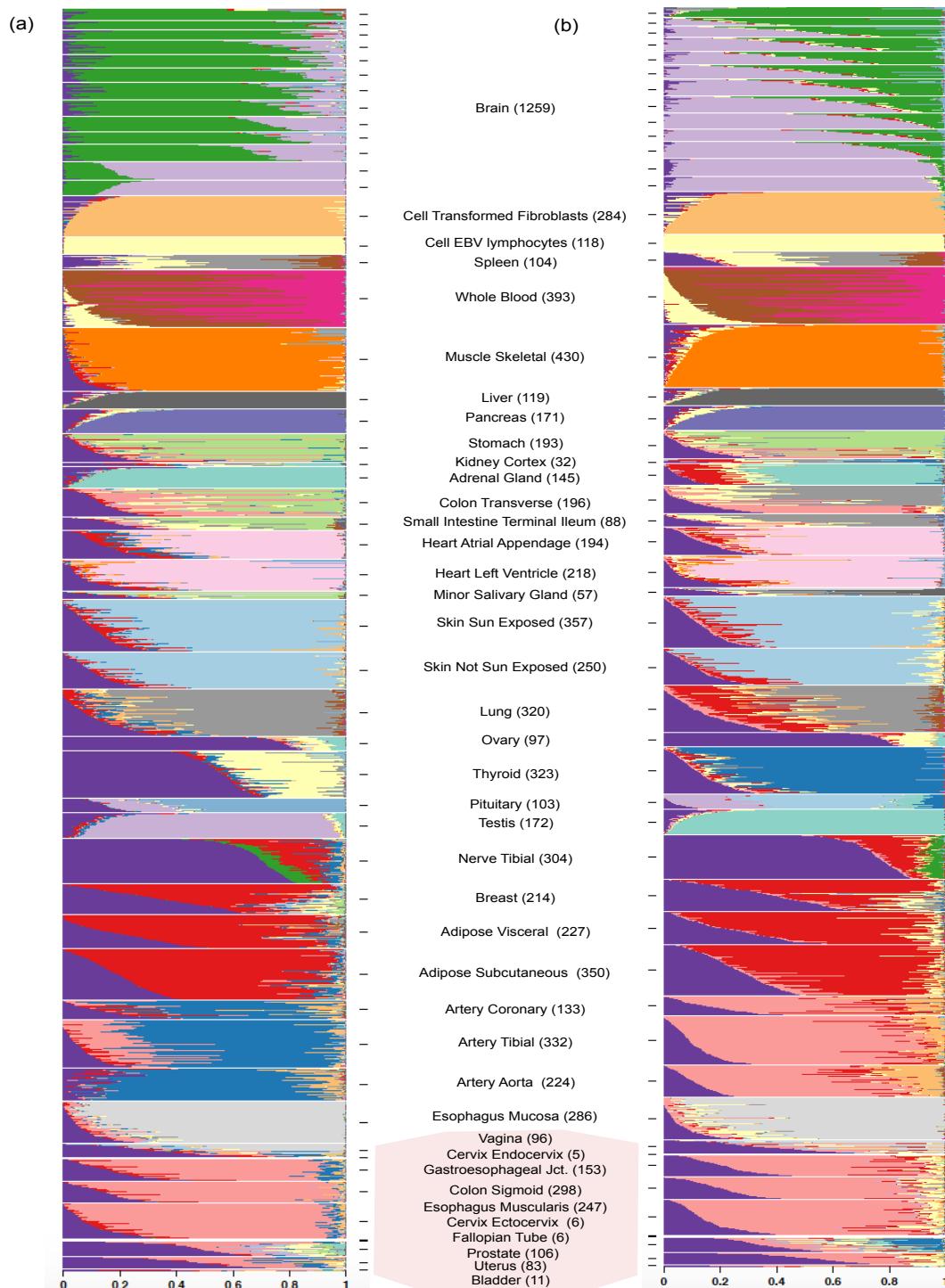


1 Supplementary figures

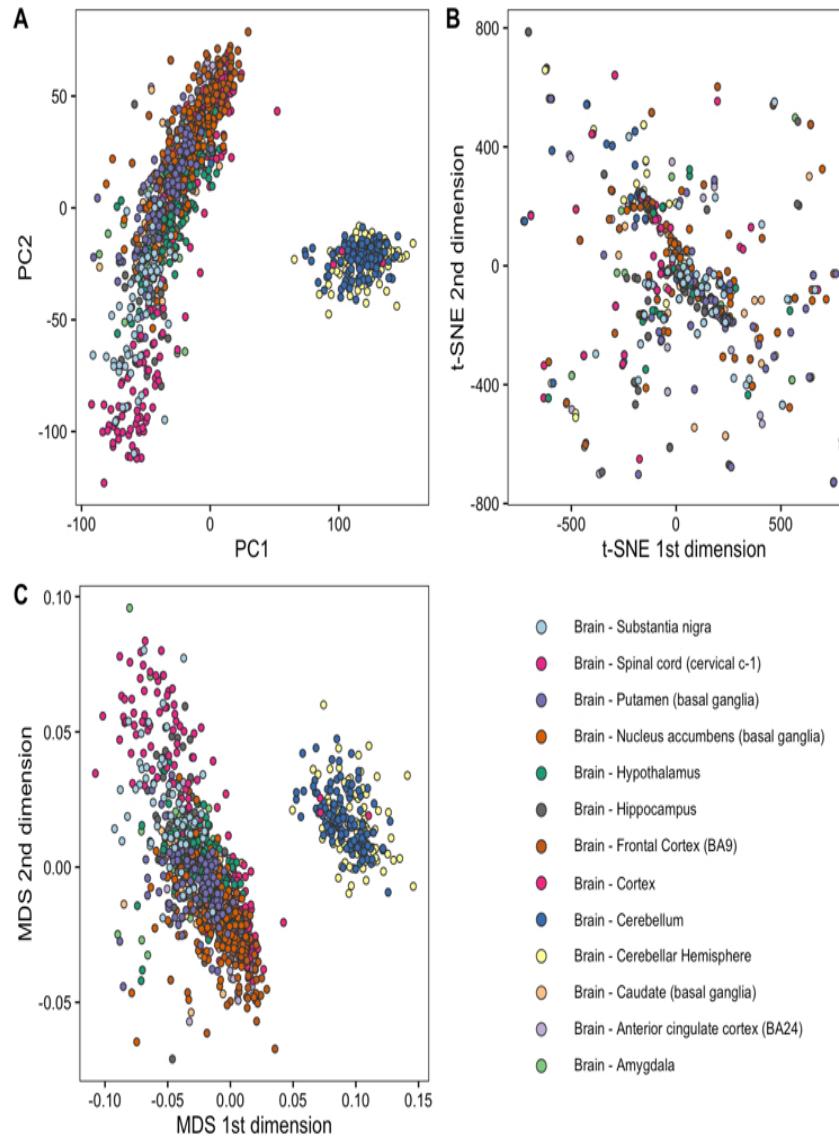
S1 Fig. Structure plot of GTEx V6 all tissue samples for (a) $K = 5$, (b) $K = 10$, (c) $K = 15$, (d) $K = 20$. Some tissues form a separate cluster from the rest of the tissues from $K = 5$ onwards (for example: Whole Blood, Skin), whereas some tissues only form a distinctive subgroup only at $K = 20$ (for example: Arteries).



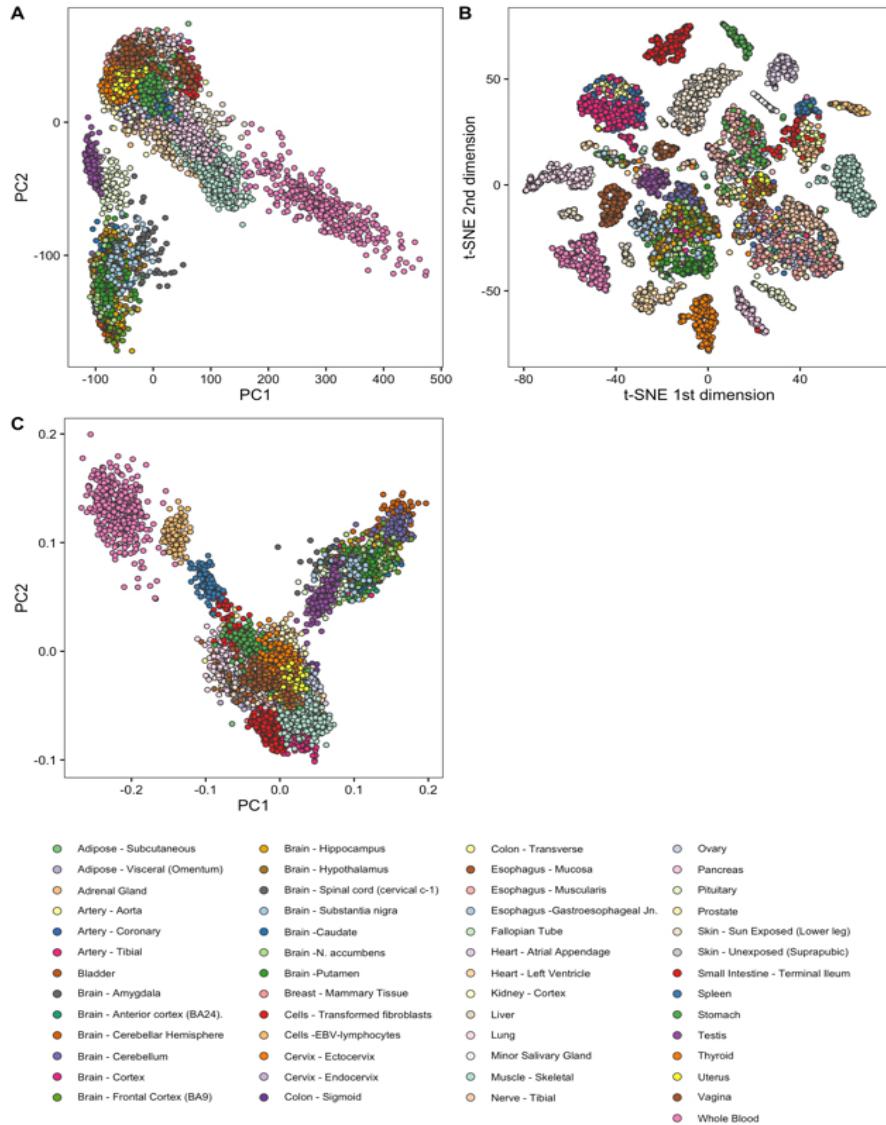
S2 Fig. Structure plot of GTEx V6 all tissue samples K=20 in 2 runs under the thinning parameters settings (a) $p_{thin} = 0.01$ and (b) $p_{thin} = 0.0001$. The patterns in two plots closely correspond to the plot in Fig 1(a), though there are a few differences from the unthinned version.



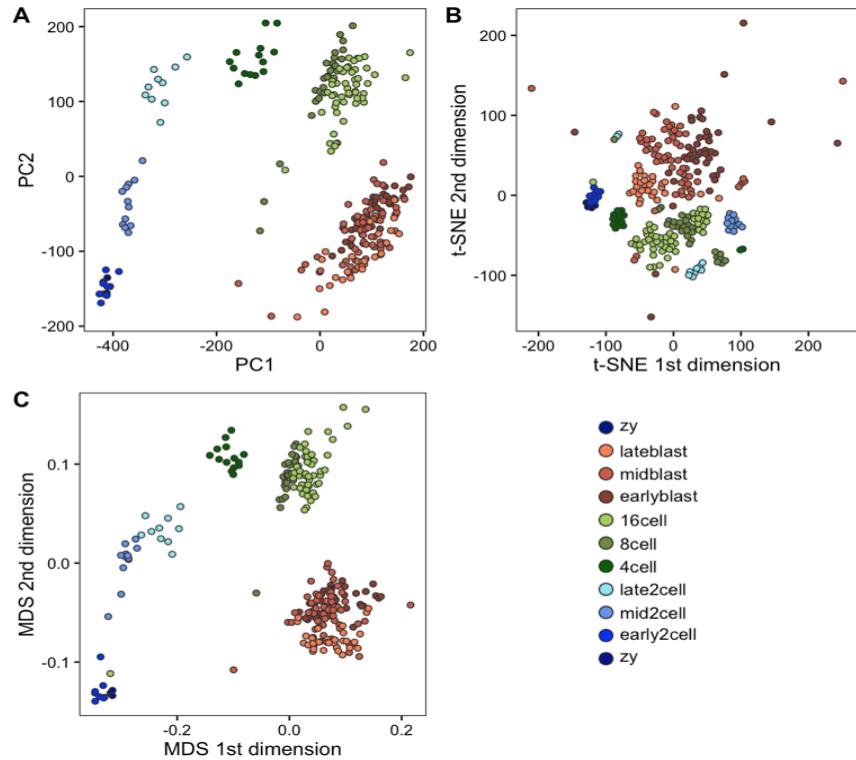
S3 Fig. GTEx brain tissues data visualization of all tissue samples using (a) principle component analysis and (b) t-SNE (c) Multidimenstional scaling. Samples of matching tissue types are indicated by points of matching color.



S4 Fig. GTEx data visualization of all tissue samples using (a) principle component analysis and (b) t-SNE (c) Multidimenstional scaling. Samples of matching tissue types are indicated by points of matching color.



S5 Fig. Mouse embryo single cell sample visualization using (a) principle component analysis and (b) t-SNE (c) Multidimensional scaling. Single cell samples collected at the same developmental stage are indicated by points of matching color.



S6 Fig. Visualization of loadings from Sparse Factor Analysis on GTEx data.

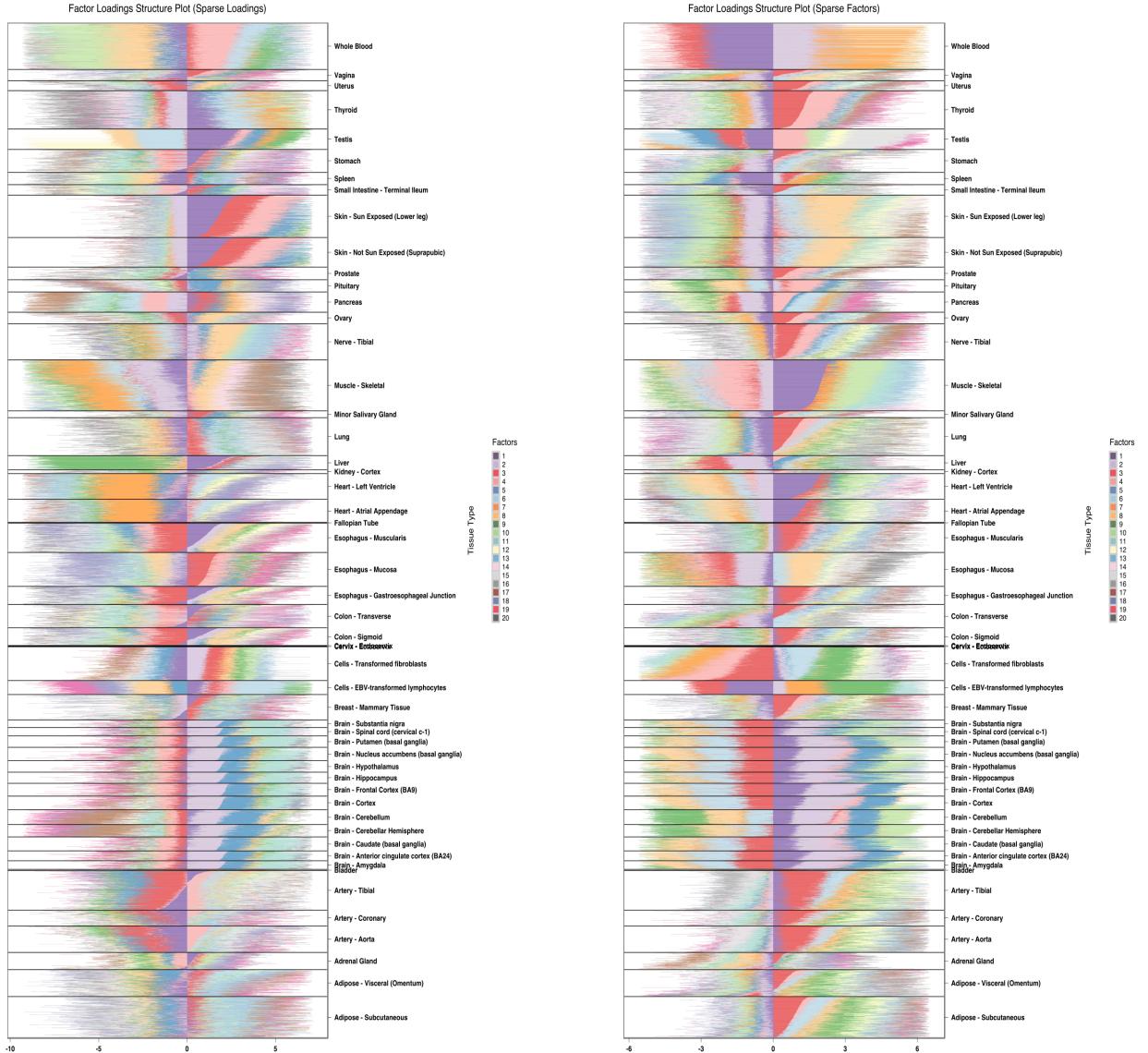


Fig 1. Stacked bar chart visualization of loadings Sparse Factor Analysis (SFA) on 8555 tissue samples across the 53 tissue types in GTEx V6 data. **(left)**: when the loadings are sparse. **(right)** when the factors are sparse.

S7 Fig. Visualization of loadings from Sparse Factor Analysis on Deng et al single cell developmental phase data.

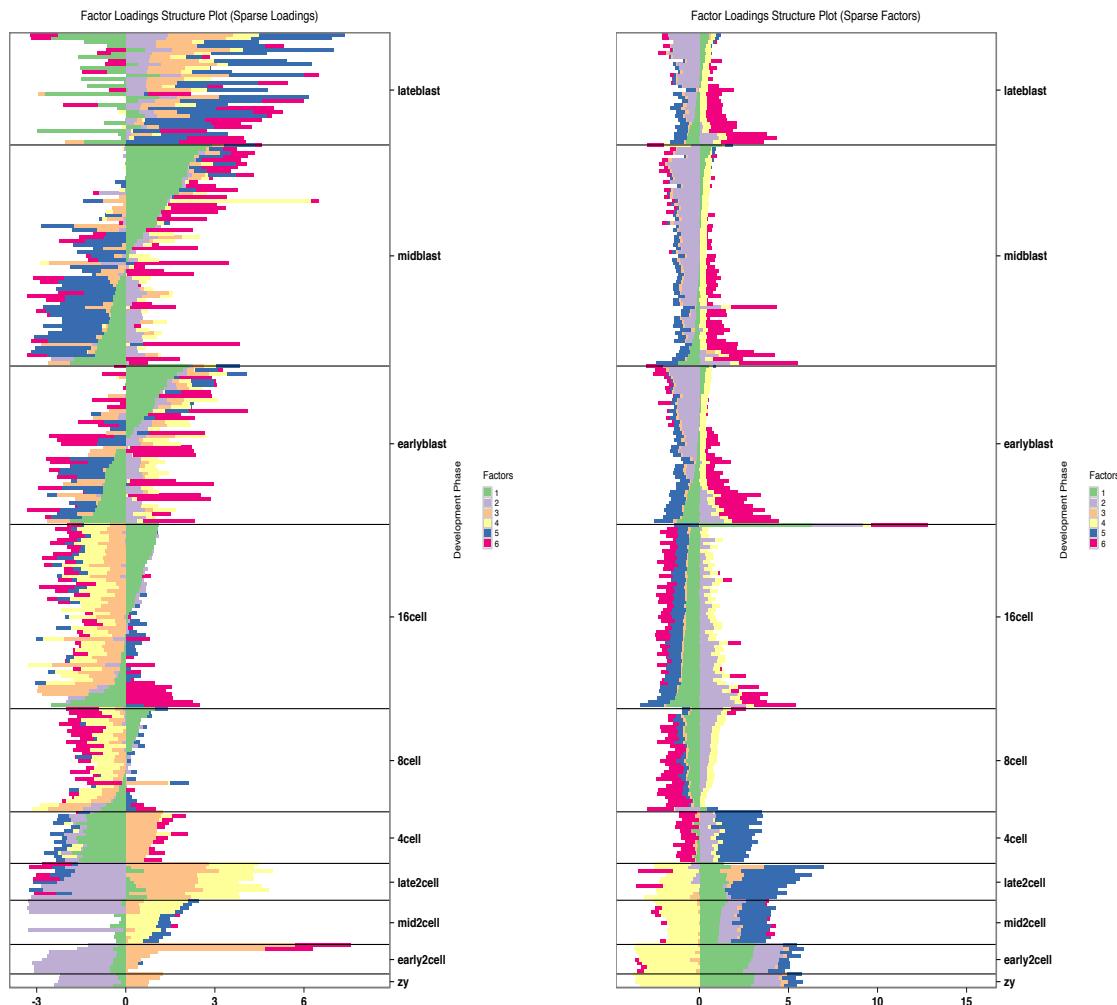


Fig 2. Stacked bar chart visualization of loadings Sparse Factor Analysis (SFA) on 259 single cell samples across developmental phases in mouse embryo for Deng *et al* data. **(left)**: when the loadings are sparse. **(right)** when the factors are sparse.

S8 Fig. Visualization of loadings from Sparse Factor Analysis on GTEx Brain data.

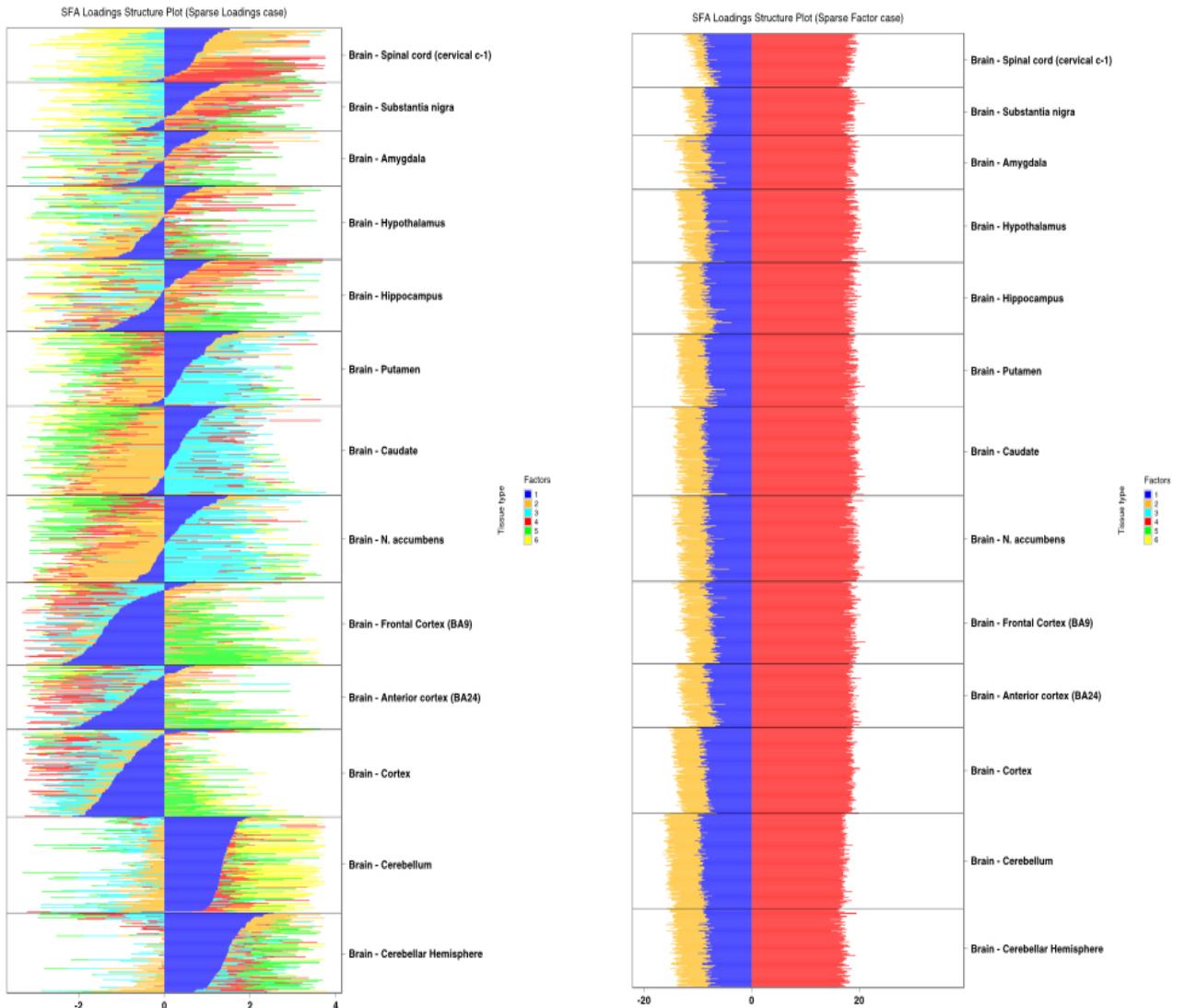


Fig 3. Stacked bar chart visualization of loadings Sparse Factor Analysis (SFA) on 1259 samples from different Brain tissues in the GTEx V6 data. **(left)**: when the loadings are sparse. **(right)** when the factors are sparse.

S9 Fig. Comparison between GoM model and hierarchical in terms of power to separate samples from pairs of tissues.

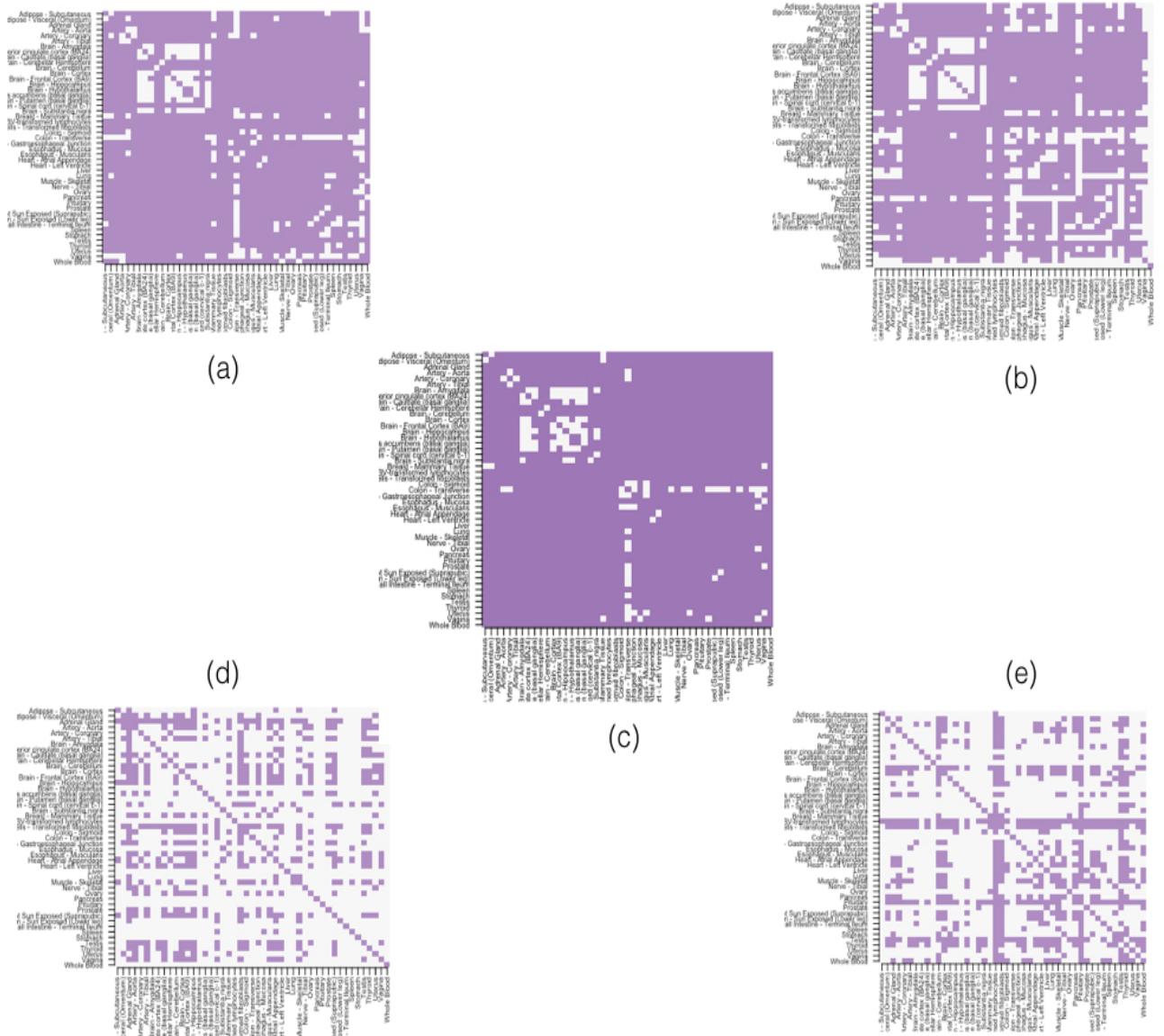


Fig 4. A comparison of accuracy of GoM model vs hierarchical clustering. Image plots to compare the GoM model with 4 different hierarchical clustering models on various transformations of the data. For each pair of tissues from the GTEx data we assessed whether or not each method (with $K = 2$ clusters) separated the samples precisely according to their actual tissue of origin, with successful separation indicated by a filled square. Very clearly, the GoM model seems to be more successful in separating pairs of tissues compared to any of the hierarchical clustering approaches. In SubFig (a), hierarchical clustering was performed on log counts per million (cpm) data using Euclidean distance. In SubFig(b), the log cpm data data was mean and scale transformed for each gene and then the hierarchical clustering was performed on the transformed data using the Euclidean distance. In SubFig (d), the hierarchical clustering was performed on counts data with the assumption the counts c_{ng} for each gene have a variance $\bar{c}_g + 1$, which we used to scale while computing distance matrix. In SubFig (e), we took the same scaled data as in SubFig(c), but we additionally performed mean and scale adjustments further so that all genes have expression of mean 0 and variance 1. In SubFig(c), GoM model is used to separate the tissues. Very clearly, GoM model seems to be performing better than any of the hierarchical methods.

S10 Fig. Dendrogram representation for Deng et al (2014) developmental phase single cell data.

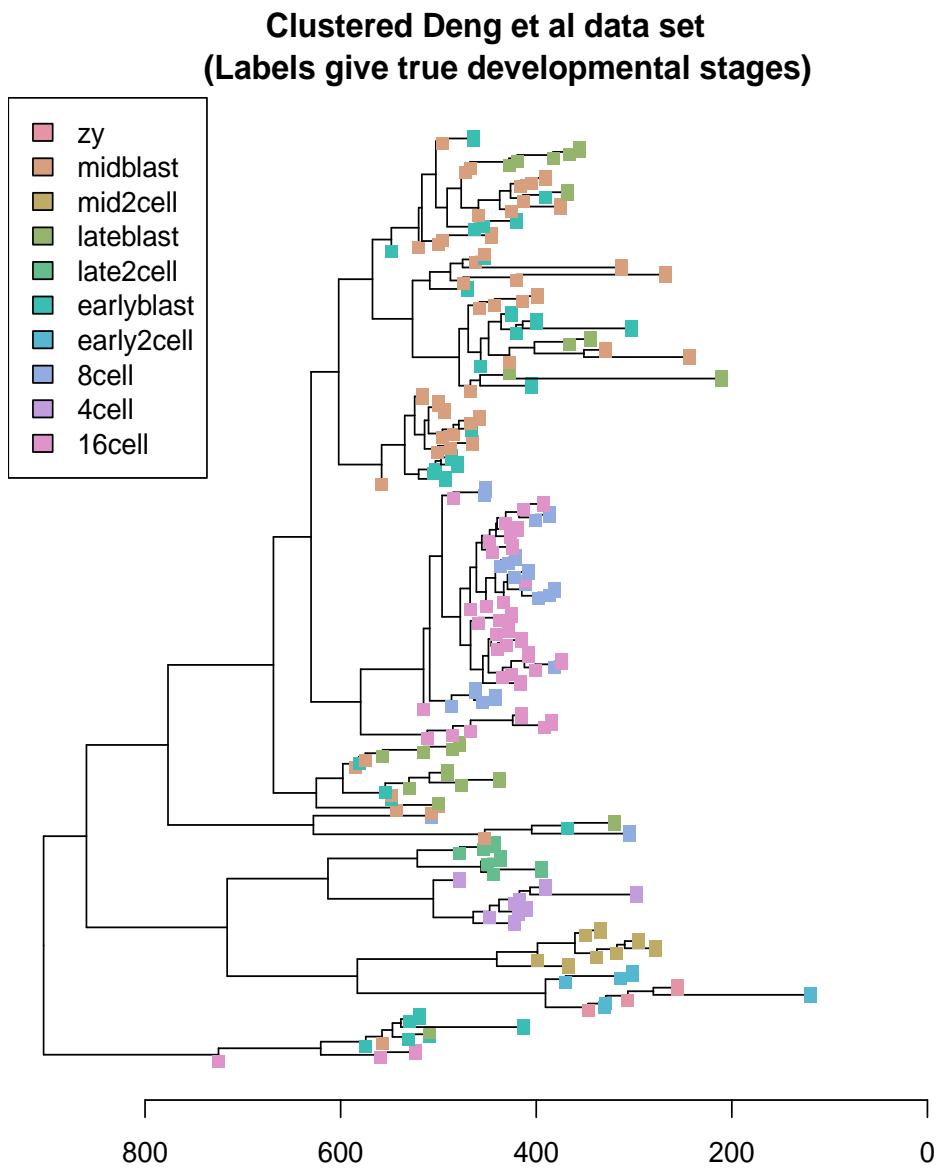


Fig 5. Dendrogram representation of the 259 cells collected across different developmental phases by Deng et al (2014). The hierarchical clustering in this case has been performed on the log counts per million (cpm) data. Note that the labels in this case are difficult to interpret and also the hierarchical clustering fails to represent the continuum among the different developmental phases as in PCA or the GoM model.

g

S11 Fig. Dendogram representation for GTEx Brain tissue level data.

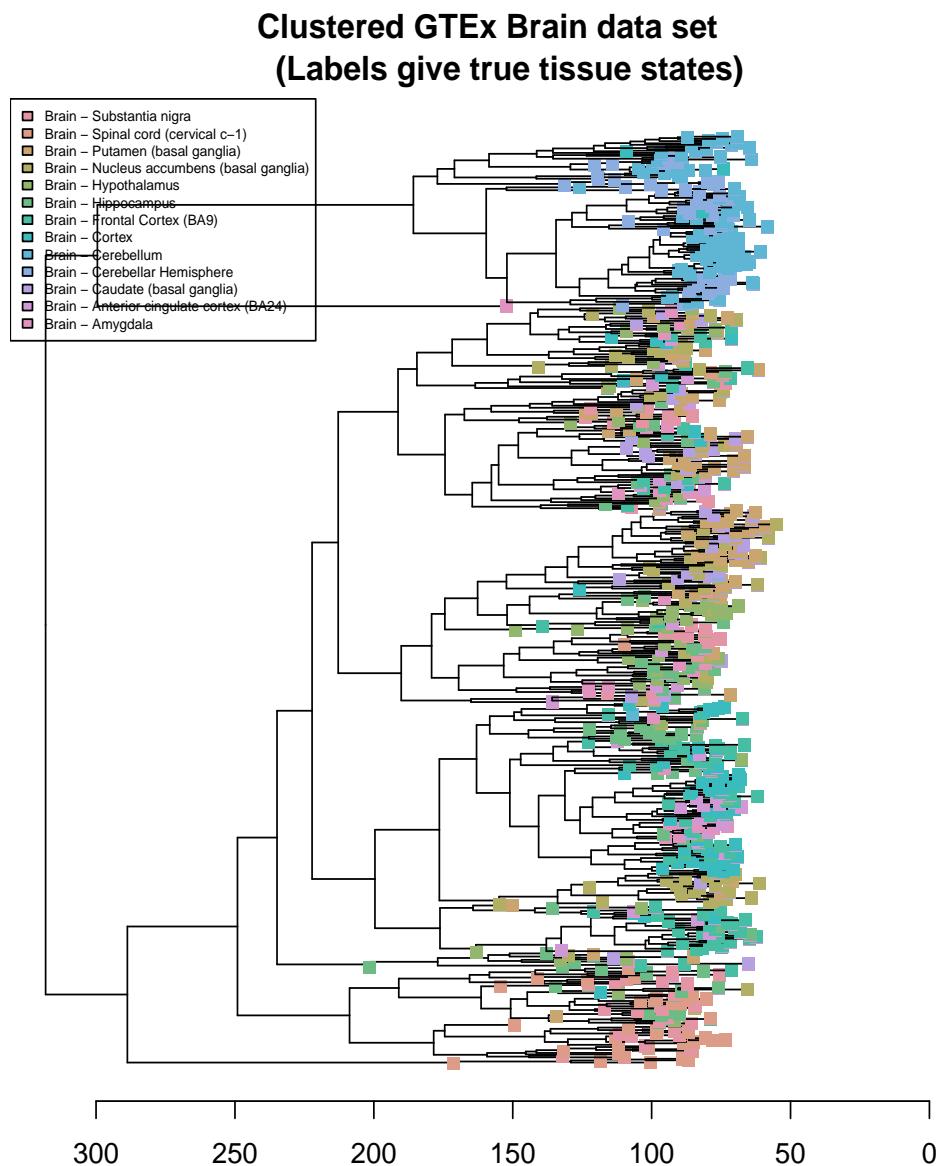


Fig 6. Dendrogram representation of the 1259 GTEx brain tissue samples. The hierarchical clustering in this case has been performed on the log counts per million (cpm) data. Note that the labels in this case are not visible because of large sample size but largely it separates out Brain Cerebellar, Cerebellar hemisphere and Brain Spinal Cord, Substantia Nigra from other parts of the brain. But any other structure in the data is not easy to see or interpret.

S12 Fig. Dendogram representation for GTEx all tissues level data.

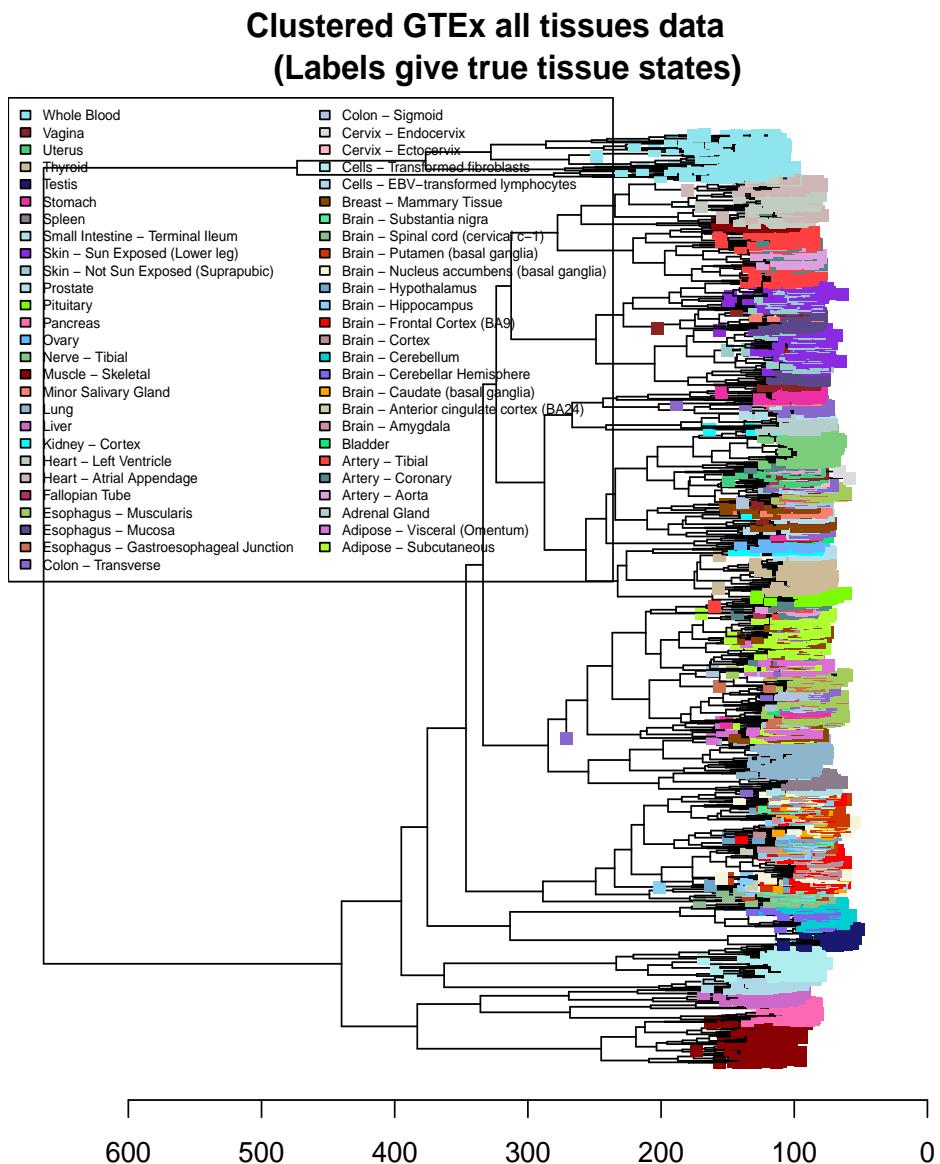


Fig 7. Dendrogram representation of the 8555 GTEx tissue samples across all the tissues. The hierarchical clustering in this case has been performed on the log counts per million (cpm) data. Very clearly, samples from different tissues seem to cluster together, but any further patterns, for instance, how similar different tissue types are, is hard to see.

S13 Fig. Top oder PC scatter plots for GTEx V6 all tissues data.

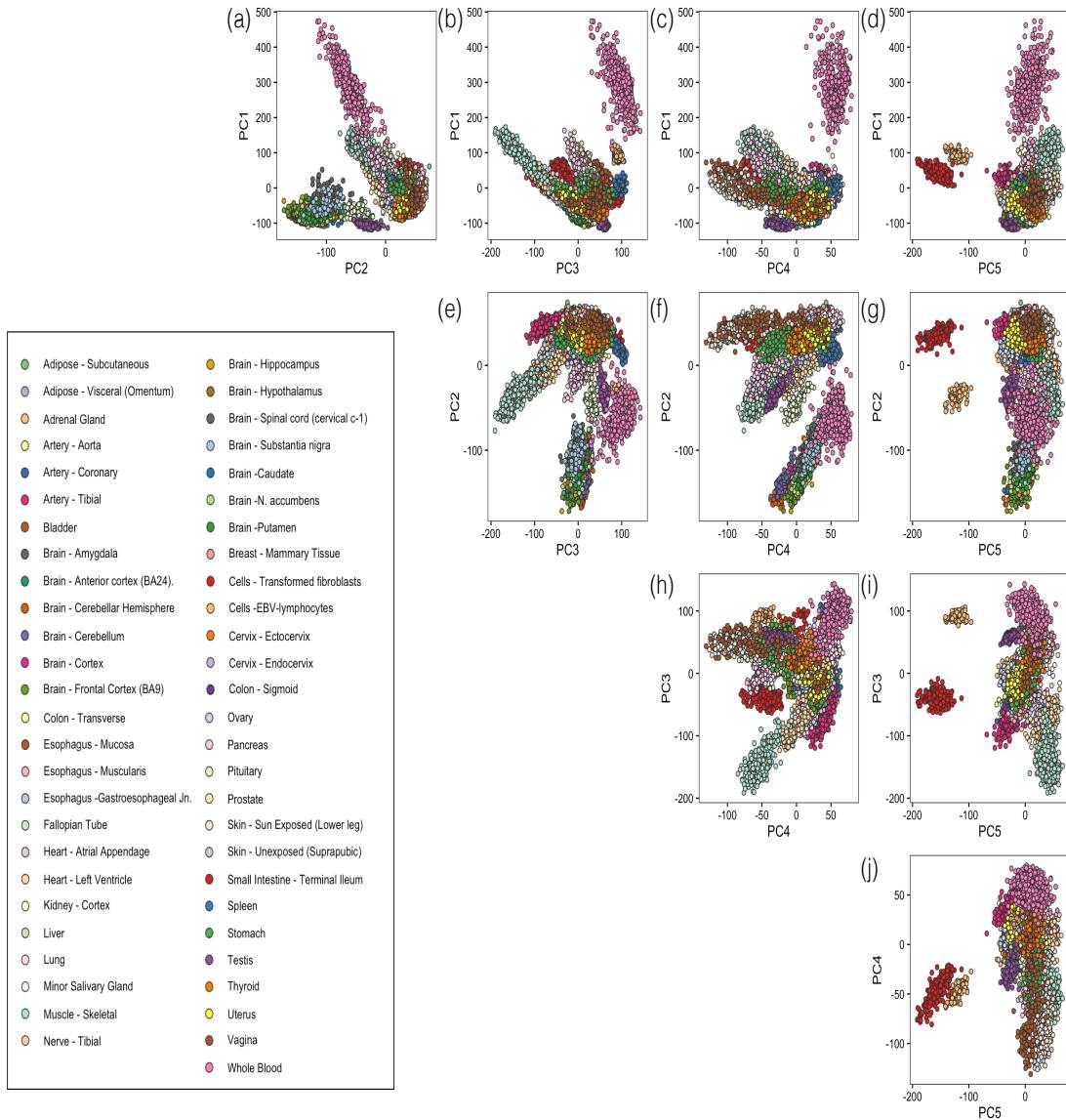
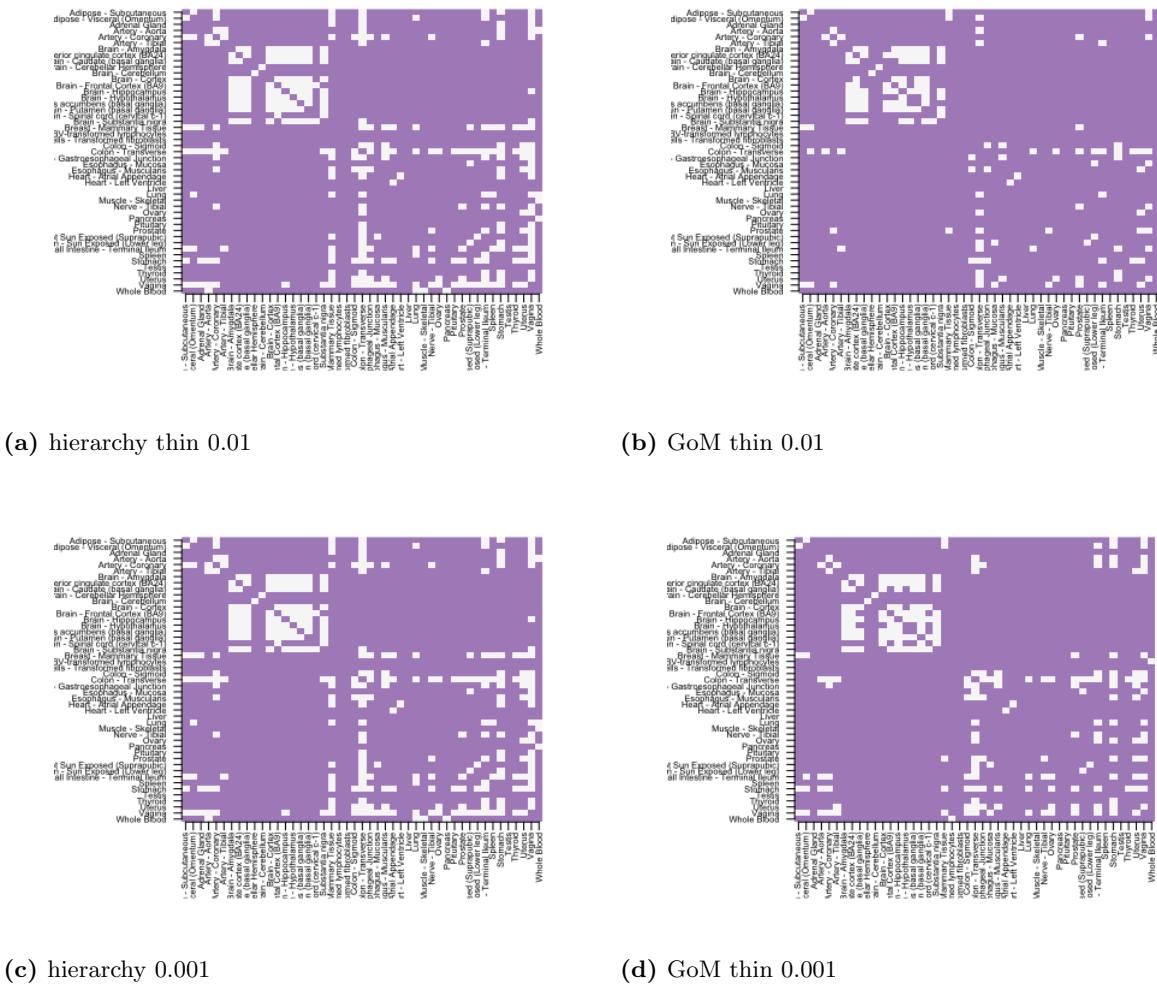


Fig 8. Scatter plot representation of the different top order PCs of the GTEx samples based on the RNA-seq expression data (log cpm normalized).

S14 Fig. A comparison of “accuracy” of hierarchical vs. model-based clustering on thinned GTEx data, with thinning parameter $p_{thin} = 0.01$ and $p_{thin} = 0.001$. For each pair of tissues from the GTEx data we assessed whether or not each clustering method (with $K = 2$ clusters) separated the samples according to their actual tissue of origin, with successful separation indicated by a filled square. Thinning deteriorates accuracy compared with the unthinned data (Fig 2), but even then the model-based method remains more successful than the hierarchical clustering in separating the samples by tissue or origin.



2 Supplementary tables

S1 Table. Cluster Annotations GTEx V6 data with top driving gene summaries.

Cluster	Top Driving Genes	Gene names	Gene Summary
1, Royal purple	<i>NEAT1</i>	nuclear paraspeckle assembly transcript 1	produces a long non-coding RNA (lncRNA) transcribed from the multiple endocrine neoplasia locus, regulates genes involved in cancer progression.
	<i>CCNL2</i>	cyclin L2	regulator of the pre-mRNA splicing process, as well as in inducing apoptosis by modulating the expression of apoptotic and antiapoptotic proteins.
	<i>SRSF5</i>	serine/arginine-rich splicing factor 5	encodes proteins of serine/arginine (SR)-rich family, involved in mRNA export from the nucleus and in translation.
2, Light purple	<i>SNAP25</i>	synaptosomal-associated protein, 25kDa	this gene product is a presynaptic plasma membrane protein involved in the regulation of neurotransmitter release.
	<i>FBXL16</i>	F-box and leucine-rich repeat protein 16	members of F-box protein family, which interact with SKP1 through the F box, and they interact with ubiquitination targets through other protein interaction domains.
	<i>SLC17A7</i>	neurochondrin	encodes proteins expressed in neuron-rich regions; associated with the membranes of synaptic vesicles and functions in glutamate transport.
3, Red	<i>FABP4</i>	fatty acid binding protein 4	encodes the fatty acid binding protein found in adipocytes, takes part in fatty acid uptake, transport, and metabolism.
	<i>PLIN1</i>	perilipin 1	protein encoded by this gene coats lipid storage droplets in adipocytes, thereby protecting them until they can be broken down by hormone-sensitive lipase.
	<i>FASN</i>	fatty acid synthase	catalyze the synthesis of palmitate from acetyl-CoA and malonyl-CoA, in the presence of NADPH, into long-chain saturated fatty acids.
4, Salmon	<i>ACTG2</i>	actin, gamma 2, smooth muscle, enteric	involved in various types of cell motility and in the maintenance of the cytoskeleton.
	<i>MYH11</i>	myosin, heavy chain 11, smooth muscle	protein encoded by this gene is a smooth muscle myosin belonging to the myosin heavy chain family, functions as a major contractile protein, converting chemical energy into mechanical energy through the hydrolysis of ATP.
	<i>SYNM</i>	synemin	protein has been found to form a linkage between desmin, which is a subunit of the IF network, and the extracellular matrix, and provides an important structural support in muscle.
5, Denim	<i>RGS5</i>	regulator of G-protein signaling 5	encodes a member of the regulators of G protein signaling (RGS) family, associated with retinal arterial macroaneurysm.
	<i>MFGE8</i>	milk fat globule-EGF factor 8 protein	encodes a preproprotein that is proteolytically processed to form multiple protein products, been implicated in wound healing, autoimmune disease, and cancer
	<i>ITGA8</i>	synemin	Proteins generated mediate numerous cellular processes including cell adhesion, cytoskeletal rearrangement, and activation of cell signaling pathways.
6, Light denim	<i>KRT10</i>	keratin 10	encodes a member of the type I (acidic) cytokeratin family, mutations associated with epidermolytic hyperkeratosis.
	<i>KRT1</i>	keratin 1, type II	specifically expressed in the spinous and granular layers of the epidermis with family member KRT10 and mutations in these genes have been associated with bullous congenital ichthyosiform erythroderma.
	<i>KRT2</i>	keratin 2, type II	expressed largely in the upper spinous layer of epidermal keratinocytes and mutations in this gene have been associated with bullous congenital ichthyosiform erythroderma.
7, Orange	<i>NEB</i>	nebulin	encodes nebulin, a giant protein component of the cytoskeletal matrix that coexists with the thick and thin filaments within the sarcomeres of skeletal muscle, associated with recessive nemaline myopathy.
	<i>MYH1</i>	myosin, heavy chain 1, skeletal muscle, adult	a major contractile protein which converts chemical energy into mechanical energy through the hydrolysis of ATP.
	<i>MYH2</i>	myosin, heavy chain 2, skeletal muscle, adult	encodes a member of the class II or conventional myosin heavy chains, and functions in skeletal muscle contraction.

