

Contents

1	Supplemental Methods S1: Modeling and Simulation	2
1.1	Simulation 1 - variance of a gene in bulk tissue	2
1.2	Simulation 2 - correlation of two genes in bulk tissue	3
2	Supplemental Methods S2: Gene variance explained by variance of the markers	5
3	Supplemental Figures	6
4	Supplemental Tables	18

1 Supplemental Methods S1: Modeling and Simulation

For a given gene q and a tissue with m cell types, we use the vector \mathbf{c}_q to represent its Cell Type expression profile (CT profile), which contains the expression level of q in each of the m cell types:

$$\mathbf{c}_q = [t_{q1}, t_{q2}, \dots, t_{qm}], \quad t_{q1}, t_{q2}, \dots, t_{qm} \geq 0, \quad \sum_{i=1}^m t_{qi} > 0 \quad (1)$$

Where t_{qi} is the mean expression level of q in cell type i . We denote the variance of values t_{qi} s by $\text{var}(t_q)$, that is, the variance of the expression level of q among different cell types.

For a given dataset with n samples from the tissue, vector α_j represents the Cell Type Composition (CTC) in sample j :

$$\alpha_j = [\alpha_{j,1}, \alpha_{j,2}, \dots, \alpha_{j,m}], \quad \alpha_{j,i} \geq 0, \quad \sum_{i=1}^m \alpha_{j,i} = 1 \quad (2)$$

Where m is the count of cell types in the tissue, as it was in definition 1 and $\alpha_{j,i}$ is the proportion of cell type i in sample j . Matrix \mathbf{A} defined as:

$$\mathbf{A} = [\alpha'_1, \alpha'_2, \dots, \alpha'_n] \quad (3)$$

contains cellular composition vectors for the n samples of the dataset.

1.1 Simulation 1 - variance of a gene in bulk tissue

Having the above definitions, the expression level of q in n samples of the dataset is given as \mathbf{e}_q :

$$\mathbf{e}_q = \mathbf{c}_q \mathbf{A} \quad (4)$$

We denote the variance of the expression level of q among the n samples of dataset presented in \mathbf{e}_q by $\text{var}(q)$. From the above formulation, it is apparent that for the two special cases when $t_{q1} = t_{q2} = \dots = t_{qm}$, (i.e $\text{var}(t_q) = 0$) or when $\alpha_1 = \dots = \alpha_n$, we have $\text{var}(q) = 0$. In a bulk tissue dataset, these two cases are equivalent to when q has the same expression level in all cell

types and when the cellular composition remains the same among the samples.

For simulation, vectors c (CT expression profiles) were generated for 1000 genes for $m = 10$ cell types. For each gene q , c_q was generated using normal distribution with mean and variance obtained from a CT expression profile of a randomly selected gene from the sNuc-seq dataset. Bulk tissue data was generated with $n = 100$ samples. Matrix A was generated using uniform distribution with the criteria that each column has sum one. The results show that among the genes, $var(t_q)$ and $var(q)$ are highly correlated (see Supplemental Fig S1).

MATLAB code for simulation

```
myVar % variance of the CT expression profiles obtained from sNuc-seq data
myMean % mean of the CT expression profiles obtained from sNuc-seq data

m = 10 % count of tissues
n = 100 % count of bulk samples

% normalizing the weight matrix
A = rand(m, n);
for i = 1:n
    A(:, i) = A(:, i) ./ (sum(A(:, i)));
end

orVar = zeros(1,1000); % variance of the CT profile
obVar = zeros(1,1000); % the observed variance in bulk tissue
for k = 1:1000 % for 1000 genes
    ind = datasample(1:length(myVar), 1);
    v = myVar(ind);
    orVar(k) = v;
    c = normrnd(myMeans(ind), sqrt(v), 1, t);
    exps = c * A;
    obVar(k) = var(exps);
end
```

1.2 Simulation 2 - correlation of two genes in bulk tissue

For given two genes p and q , with expression vectors e_p and e_q in the bulk tissue dataset, the higher the correlation of their CT profiles c_p and c_q (in any direction), the more likely that e_p and e_q are also correlated in the same direction. Having centered c_p and c_q as $c_p^c = c_p - \bar{c}_p$ and $c_q^c = c_q - \bar{c}_q$ and formulation (4), since $corr(c_p, c_q) = corr(c_q^c, c_p^c) = \cos(c_p^c, c_q^c)$; correlation

between \mathbf{e}_p and \mathbf{e}_q can be written as:

$$\begin{aligned} \text{corr}(\mathbf{e}_p, \mathbf{e}_q) &= \text{corr}([\mathbf{c}_p^c \cdot \alpha'_1, \mathbf{c}_p^c \cdot \alpha'_2, \dots, \mathbf{c}_p^c \cdot \alpha'_n], [\mathbf{c}_q^c \cdot \alpha'_1, \mathbf{c}_q^c \cdot \alpha'_2, \dots, \mathbf{c}_q^c \cdot \alpha'_n]) \\ &= \text{corr}([\|\alpha_1\| \cos(\mathbf{c}_p^c, \alpha_1), \|\alpha_2\| \cos(\mathbf{c}_p^c, \alpha_2), \dots, \|\alpha_n\| \cos(\mathbf{c}_p^c, \alpha_n)] \quad (5) \\ &\quad [\|\alpha_1\| \cos(\mathbf{c}_q^c, \alpha_1), \|\alpha_2\| \cos(\mathbf{c}_q^c, \alpha_2), \dots, \|\alpha_n\| \cos(\mathbf{c}_q^c, \alpha_n)]) \end{aligned}$$

Where $\cos(a, b)$ is the cosine of the angle between the two vectors. As the angle between \mathbf{c}_p^c and \mathbf{c}_q^c gets smaller (equivalent to their correlation getting higher), the difference between $\cos(\mathbf{c}_p^c, \alpha_i)$ and $\cos(\mathbf{c}_q^c, \alpha_i)$ gets smaller and therefore correlation between \mathbf{e}_p and \mathbf{e}_q gets higher.

Figure 3C shows simulation results for correlation of CT expression profiles versus the observed correlation in the bulk tissue for 1000 gene pairs.

MATLAB code for the simulation

```
n = 100
CTCorr = zeros(1, 1000); % correlation of the CT expression profiles
sampleCorr = zeros(1, 1000); % correlation of the gene pair in the bulk tissue
for k = 1:1000
    P = (randn(1,10)*3 + 4); % generating CT profile for gene P
    Q = (randn(1,10)*3 + 4); % generating CT profile for gene Q
    Gs = [P, Q];

    CTCorr(k) = corr(P', Q');

    A = rand(10, n);
    % normalizing the weight matrix
    for i = 1:n
        A(:, i) = A(:, i) ./ (sum(A(:,i)));
    end

    % getting the final vectors
    exps = Gs * A;

    sib = corr(exps');
    sampleCorr(k) = sib(1, 2);
end
```

2 Supplemental Methods S2: Gene variance explained by variance of the markers

We examined the results from the Principal Component Regression method from the two sets of marker genes from the same five cell types - Astrocyte, Microglia, Oligodendrocyte, Endothelial and Pyramidal. This was to study the specificity of the model in capturing cellular composition induced variance (reflected specifically in the variation of the marker genes) as oppose to variance induced by a general confounding factor shared among all the genes. The two sets of marker genes are identified independently and have a small overlap (see Supplemental Fig. S3).

3 Supplemental Figures

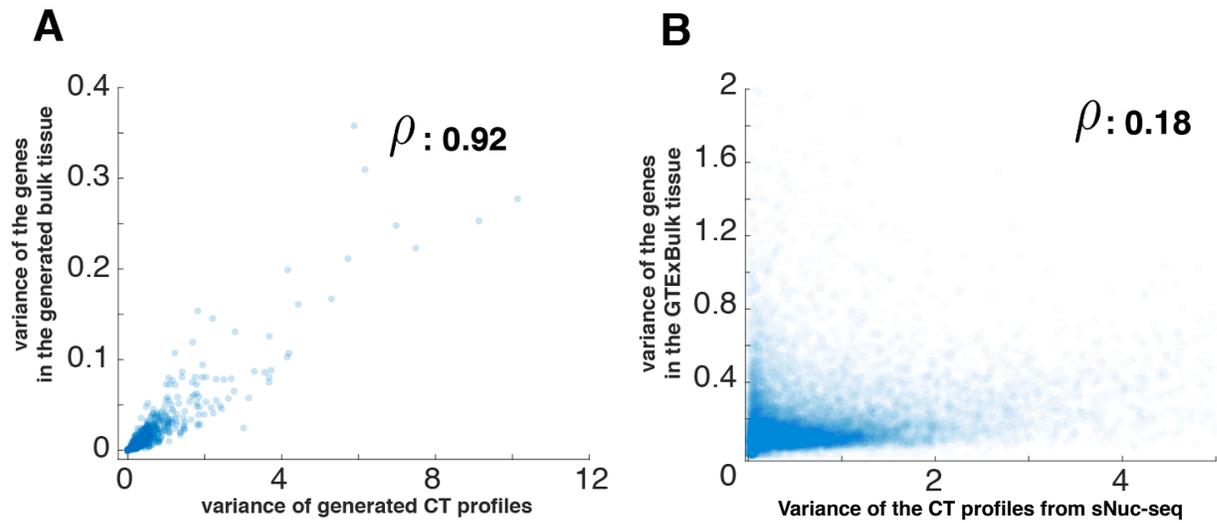


Figure S1: Variance of the CT profiles versus the variance observed in the bulk tissue. (A) Results from simulation for 1000 genes. Each point is data from a gene. **(B)** Results from data: CT profiles estimated using sNuc-seq data, compared to the observed variance in the GTExBulk tissue dataset.

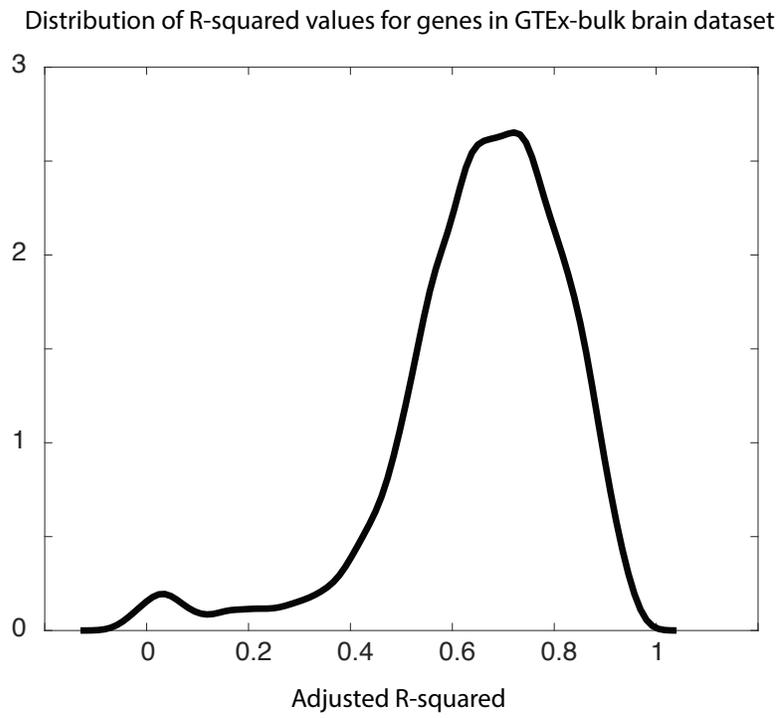


Figure S2: Distribution of the R-squared values for genes in the GTExBulk coexpression network.

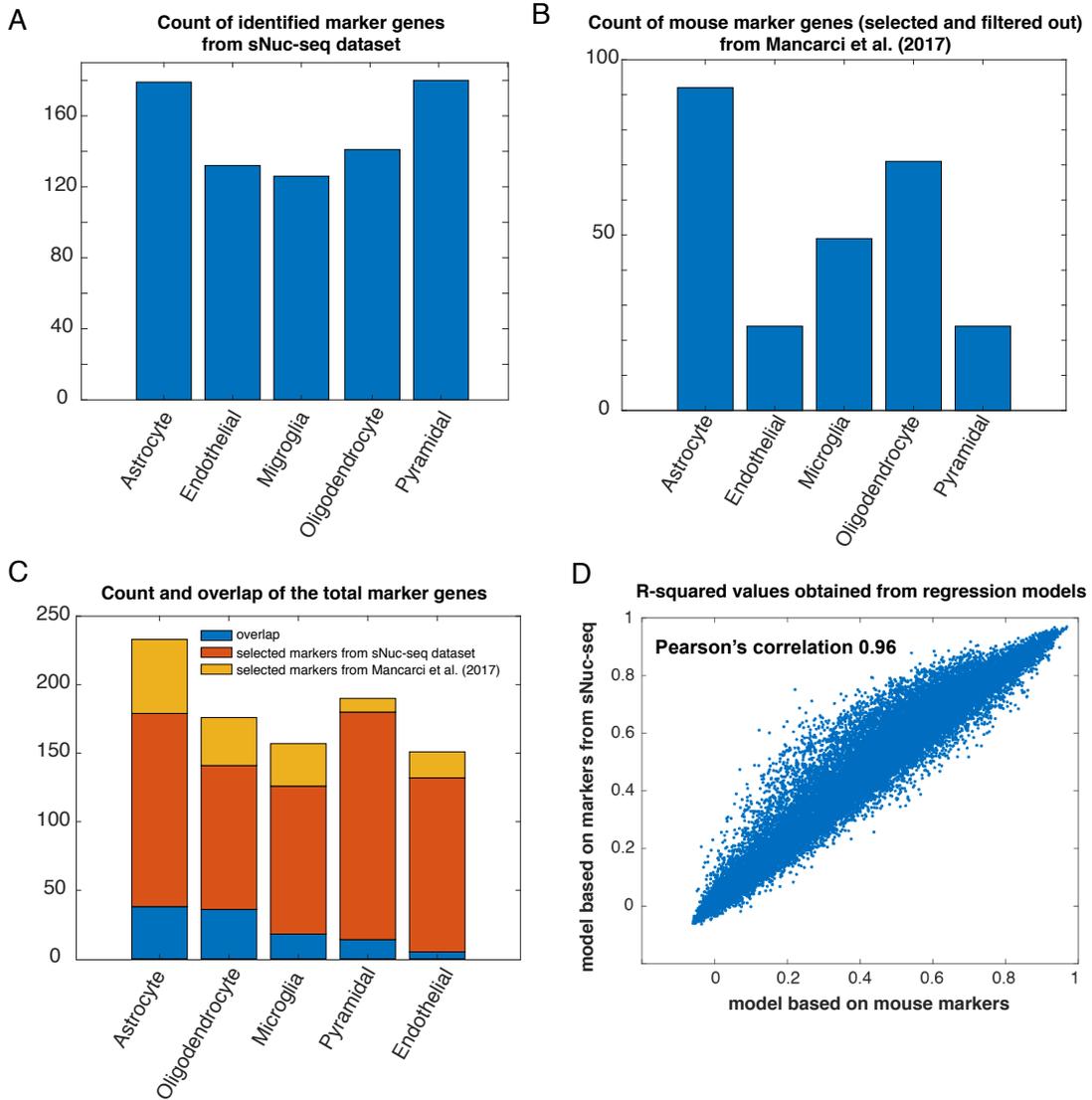


Figure S3: Comparison of two sets of marker genes. (A) Count of marker genes identified in sNuc-seq dataset. (B) Marker genes from Mancarci et al. (2017). (C) Overlap of the marker genes between the two sources. (D) R-squared values from the two sets of marker genes are highly correlated.

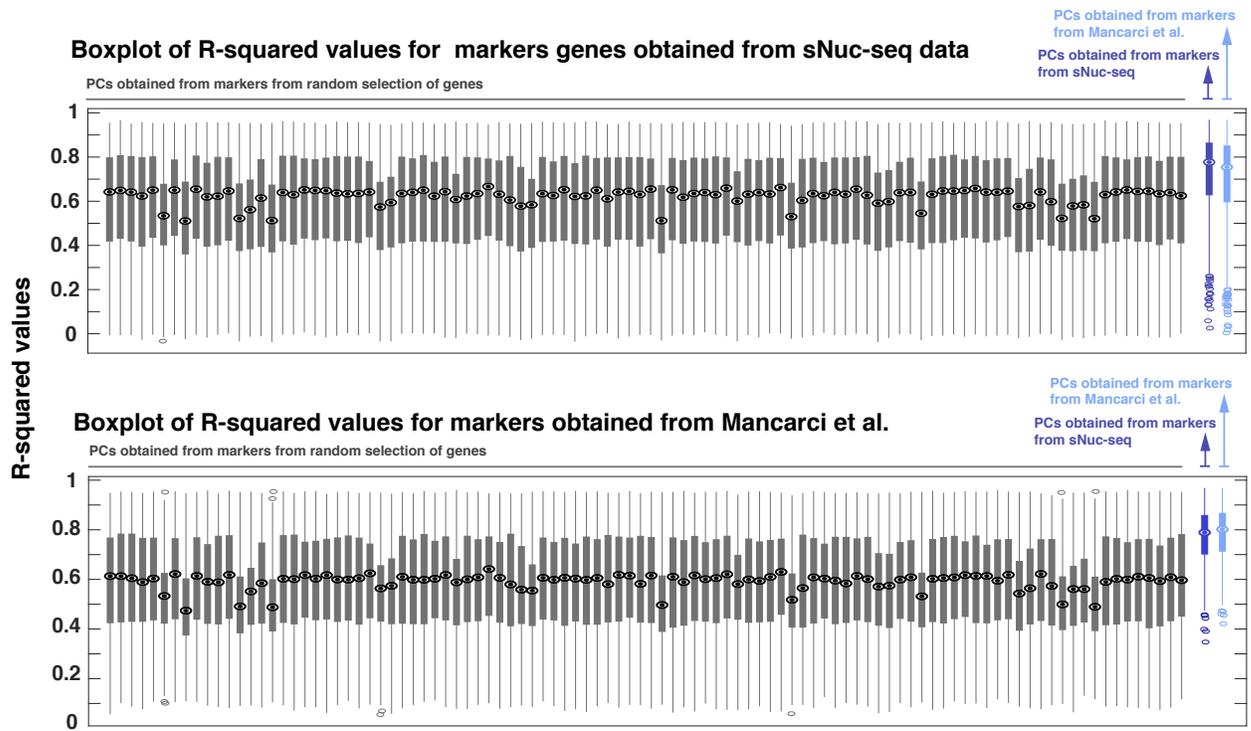


Figure S4: Comparison of the distribution of R^2 values for marker genes, when PCs are obtained from random selection of genes versus when they are obtained from marker genes.

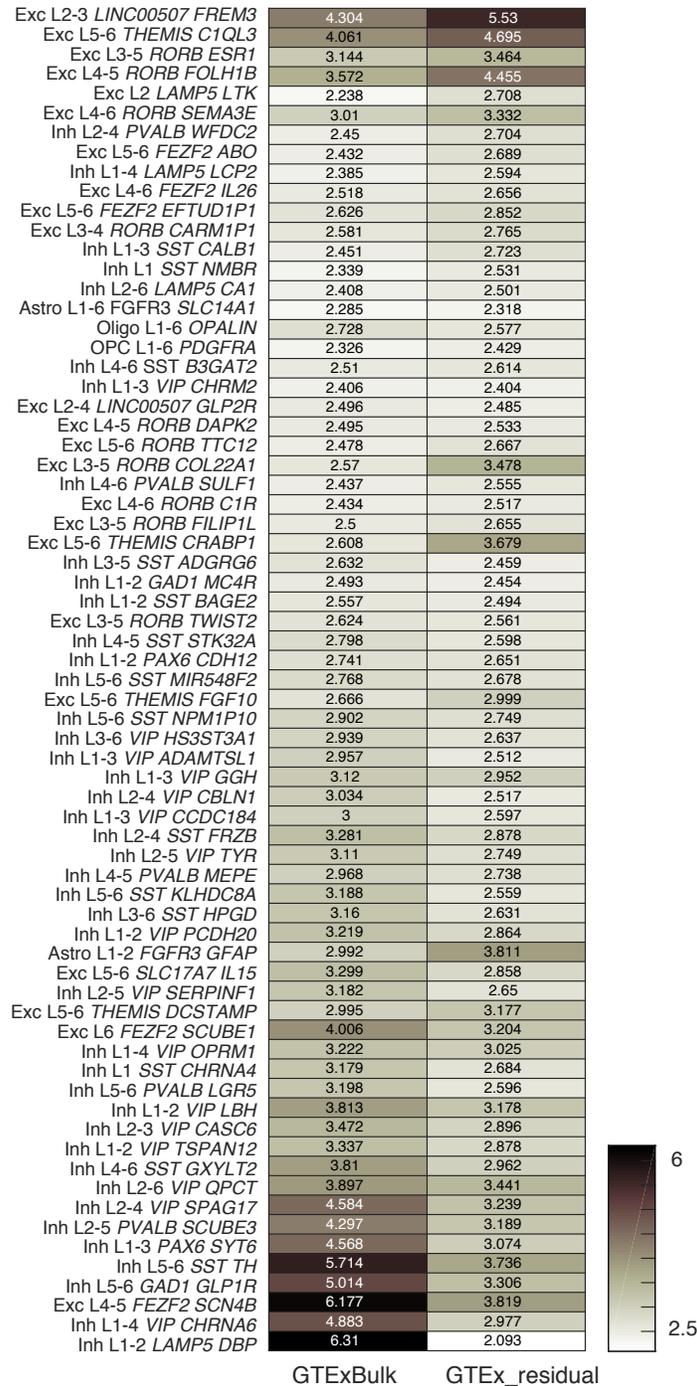


Figure S6: Ratio of the observed versus the expected link overlap of the links from sNuc-seq networks with GTExBulk and GTEX_residual networks. sNuc-seq populations are sorted based on the count of cells, with the top one having the highest count of cells. Generally, the link overlap between GTExBulk and GTEX_residual doesn't change much and there is some level of agreement between sNuc-seq population networks and the two bulk networks.

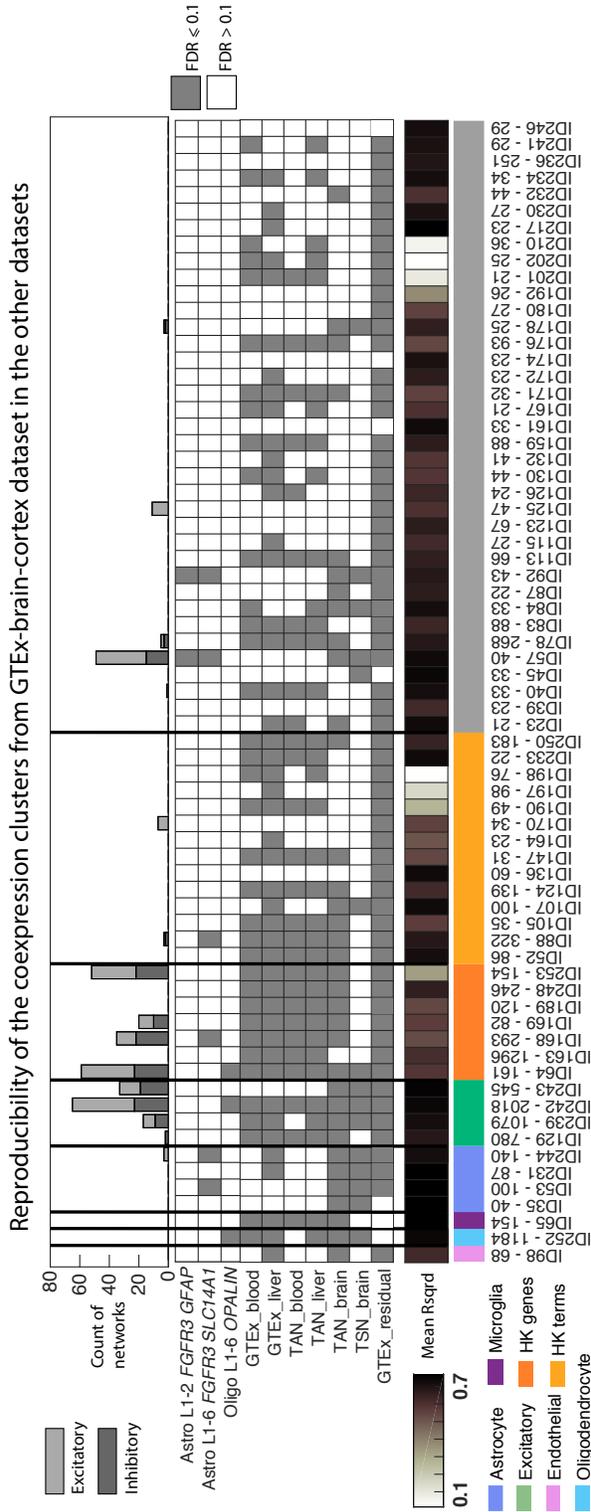


Figure S7: Reproducibility of clusters identified in GTExBulk in other networks. Rows in the heatmap show the reproducibility of clusters in different networks. Each column in the heatmap (and the same column in the bar plot) shows data for one cluster identified in GTExBulk in network. The gray color in the heatmap shows that the cluster has significantly high count of links (FDR ≤ 0.1). The top bar plot shows, for each cluster, the count of Excitatory and Inhibitory networks (built from populations of cell types identified in sNuc-seq data) in which the cluster has significantly high count of links (FDR ≤ 0.1). The bottom color bar shows if the cluster has markers of specific cell types, enriched by housekeeping genes or functional terms.

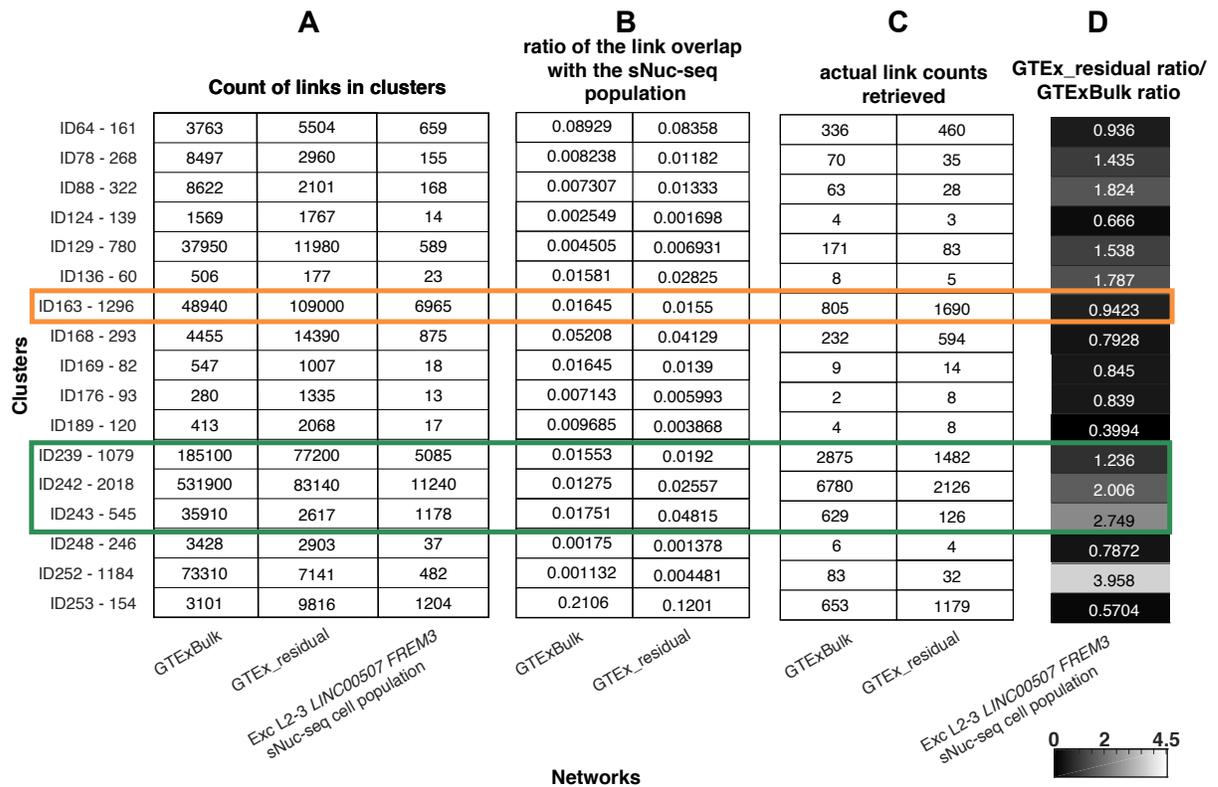


Figure S8: Representation of the links in a sNuc-seq population (Exc L2-3 *LINC00507 FREM3*). This population was selected because it has the highest count of cells among all the sNuc-seq populations and therefore a more complete network regarding the gene count. **A** shows count of links in each of the GTEXBulk clusters (clusters were selected to have 10 or more links in the sNuc-seq population). **B** shows ratio of links from GTEXBulk and GTEX_Residual networks that overlap with links in the sNuc-seq network. **C** shows the actual count of links from sNuc-seq network retrieved (has overlap) in GTEXBulk and GTEX_residual networks. **D** shows the proportion of the ratios in **B**. For the Pyramidal clusters in the green box, although the count of links has decreased, precision has almost doubled for the GTEX_residual network. For the housekeeping cluster in the orange box, count of links retrieved is more than double in GTEX_residual network compared to GTEXBulk network and the precision has not changed much.

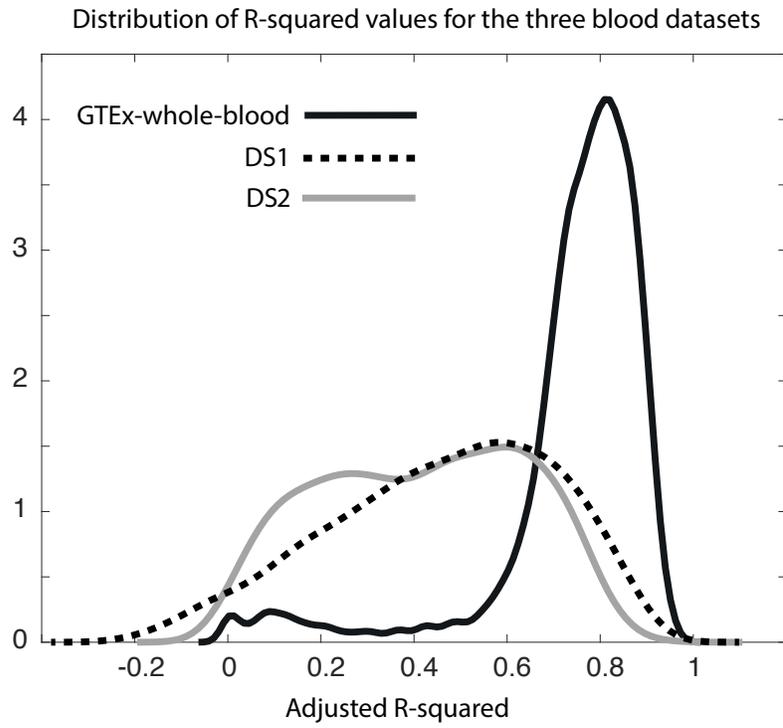


Figure S9: Distribution of the R-squared values for genes in three blood datasets.

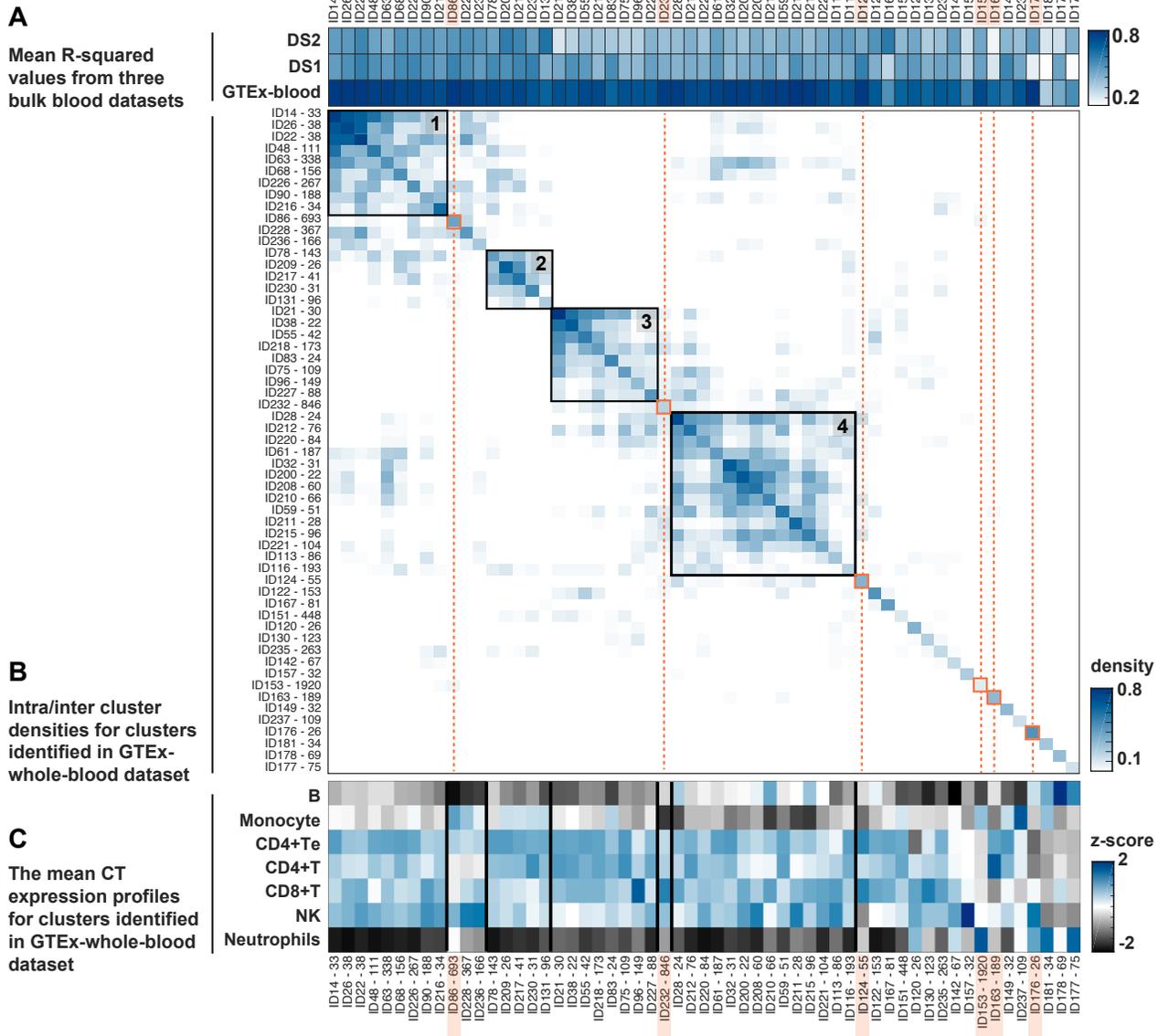
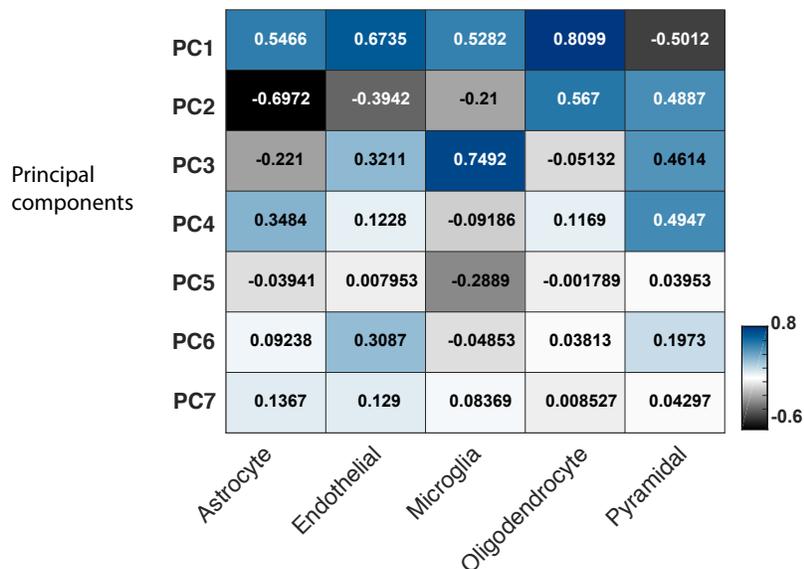


Figure S10: Cell type expression profiles of genes in whole blood dataset can explain grouping of the genes in coexpression clusters. (A) The heatmap shows the average R^2 values for coexpression clusters (identified in GTEx-whole-blood network) in GTEx-whole-blood and two microarray bulk blood datasets (DS1 and DS2). Clusters are labeled with their ID and count of genes. Clusters with > 2 genes identified as differentially expressed among multiple blood cell types are highlighted with color orange. The dotted lines trace their inter and intra cluster density in the heatmap in B. (B) The heatmap shows grouping of coexpression clusters identified in GTEx-whole-blood dataset. Four major groups of clusters are identified and labeled. (C) The heatmap shows zscores for average CT expression profiles for the clusters. All the four groups have low expression levels in Neutrophils, while expression level of Monocytes and B cells vary between them.

A

Correlation of the mean expression level of different sets of marker genes with the first seven principal component scores

**B**

Percent of variation explained by each of the principal components in the bulk tissue dataset for different sets of marker genes

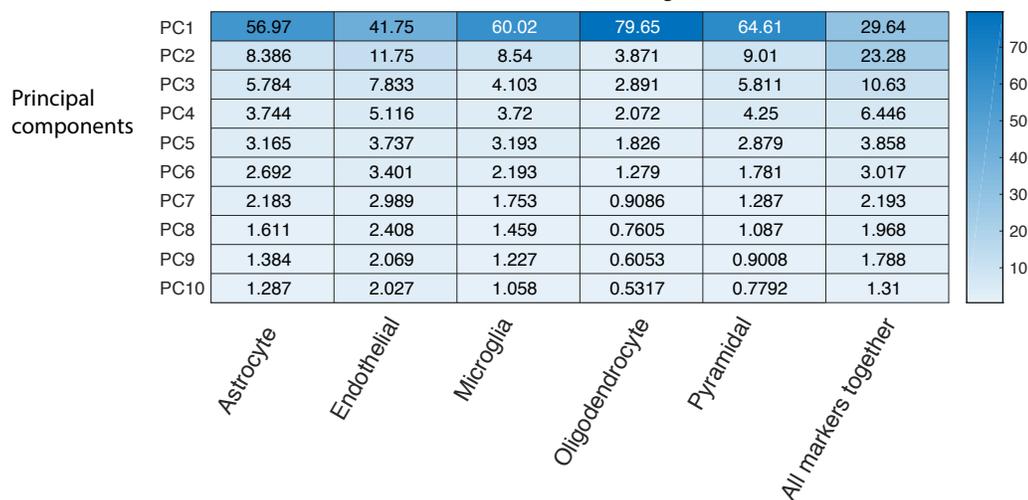


Figure S12: Cellular compositional variance captured by principal components. **A** shows correlation of the mean zscore of expression of sets of marker genes among the GTExBulk (brain) samples, with the first seven principal component scores obtained from the whole set of marker genes. **B** shows the percent of variation explained by first ten principal components of different marker sets as well as all markers together.

4 Supplemental Tables

Table S1: Count of links and genes in each of the networks

networks	gene count	link count
GTEX brain cortex	14,102	1,774,291
GTEX blood	11,348	1,312,228
GTEX liver	16,388	1,481,028
TAN brain	8,761	1,118,791
TAN blood	8,747	380,875
TAN liver	10,227	496,476
TSN brain	6,422	358,531

Exc L2-3 <i>LINC00507 FREM3</i>	10610	563000
Exc L5-6 <i>THEMIS C1QL3</i>	10710	573300
Exc L3-5 <i>RORB ESR1</i>	8509	362000
Exc L4-5 <i>RORB FOLH1B</i>	11040	609700
Exc L2 <i>LAMP5 LTK</i>	8530	363800
Exc L4-6 <i>RORB SEMA3E</i>	9194	422600
Inh L2-4 <i>PVALB WFDC2</i>	9107	410900
Exc L5-6 <i>FEZF2 ABO</i>	10140	508900
Inh L1-4 <i>LAMP5 LCP2</i>	8169	326000
Exc L4-6 <i>FEZF2 IL26</i>	9410	434500
Exc L5-6 <i>FEZF2 EFTUD1P1</i>	10810	573000
Exc L3-4 <i>RORB CARM1P1</i>	10340	515300
Inh L1-3 <i>SST CALB1</i>	5968	157300
Inh L1 <i>SST NMBR</i>	7839	278400
Inh L2-6 <i>LAMP5 CA1</i>	8656	342100
Astro L1-6 <i>FGFR3 SLC14A1</i>	4427	70560
Oligo L1-6 <i>OPALIN</i>	3286	32340
OPC L1-6 <i>PDGFRA</i>	3281	29850
Inh L4-6 <i>SST B3GAT2</i>	7402	198500
Inh L1-3 <i>VIP CHR12</i>	7583	201700
Exc L2-4 <i>LINC00507 GLP2R</i>	8894	317200
Exc L4-5 <i>RORB DAPK2</i>	8636	285700
Exc L5-6 <i>RORB TTC12</i>	10480	464400
Exc L3-5 <i>RORB COL22A1</i>	8420	252100
Inh L4-6 <i>PVALB SULF1</i>	8339	246800
Exc L4-6 <i>RORB C1R</i>	9216	323000
Exc L3-5 <i>RORB FILIP1L</i>	9276	321500
Exc L5-6 <i>THEMIS CRABP1</i>	10150	411200
Inh L3-5 <i>SST ADGRG6</i>	7035	130400
Inh L1-2 <i>GAD1 MC4R</i>	6328	79800
Inh L1-2 <i>SST BAGE2</i>	5777	64540
Exc L3-5 <i>RORB TWIST2</i>	7885	147100
Inh L4-5 <i>SST STK32A</i>	5401	55640
Inh L1-2 <i>PAX6 CDH12</i>	7204	118000
Inh L5-6 <i>SST MIR548F2</i>	8185	177600
Exc L5-6 <i>THEMIS FGF10</i>	9303	271200
Inh L5-6 <i>SST NPM1P10</i>	7086	116700
Inh L3-6 <i>VIP HS3ST3A1</i>	6367	84610
Inh L1-3 <i>VIP ADAMTSL1</i>	5939	73990
Inh L1-3 <i>VIP GGH</i>	4470	42500
Inh L2-4 <i>VIP CBLN1</i>	3773	29240
Inh L1-3 <i>VIP CCDC184</i>	5904	78420
Inh L2-4 <i>SST FRZB</i>	4674	48690
Inh L2-5 <i>VIP TYR</i>	4724	48510
Inh L4-5 <i>PVALB MEPE</i>	7049	127800
Inh L5-6 <i>SST KLHDC8A</i>	5764	74940
Inh L3-6 <i>SST HPGD</i>	5986	87560
Inh L1-2 <i>VIP PCDH20</i>	4353	42050
Astro L1-2 <i>FGFR3 GFAP</i>	917	2964
Exc L5-6 <i>SLC17A7 IL15</i>	6698	119400
Inh L2-5 <i>VIP SERPINF1</i>	2972	20940
Exc L5-6 <i>THEMIS DCSTAMP</i>	7535	173700
Exc L6 <i>FEZF2 SCUBE1</i>	6841	138800
Inh L1-4 <i>VIP OPRM1</i>	2960	20860
Inh L1 <i>SST CHR14</i>	3136	23040
Inh L5-6 <i>PVALB LGR5</i>	3522	27870
Inh L1-2 <i>VIP LBH</i>	4474	51470
Inh L2-3 <i>VIP CASC6</i>	4596	54250
Inh L1-2 <i>VIP TSPAN12</i>	2550	16730
Inh L4-6 <i>SST GXYLT2</i>	3315	29010
Inh L2-6 <i>VIP QPCT</i>	1717	9126
Inh L2-4 <i>VIP SPAG17</i>	1074	4576
Inh L2-5 <i>PVALB SCUBE3</i>	1516	7848
Inh L1-3 <i>PAX6 SYT6</i>	2202	15420
Inh L5-6 <i>SST TH</i>	1450	8536
Inh L5-6 <i>GAD1 GLP1R</i>	856	3302
Exc L4-5 <i>FEZF2 SCN4B</i>	2936	32310
Inh L1-4 <i>VIP CHR14</i>	409	1300
Inh L1-2 <i>LAMP5 DBP</i>	158	304

Count of genes
in the networks

Count of links
in the networks

Table S2: Count of genes (with one or more links) and links in different sNuc-seq networks.