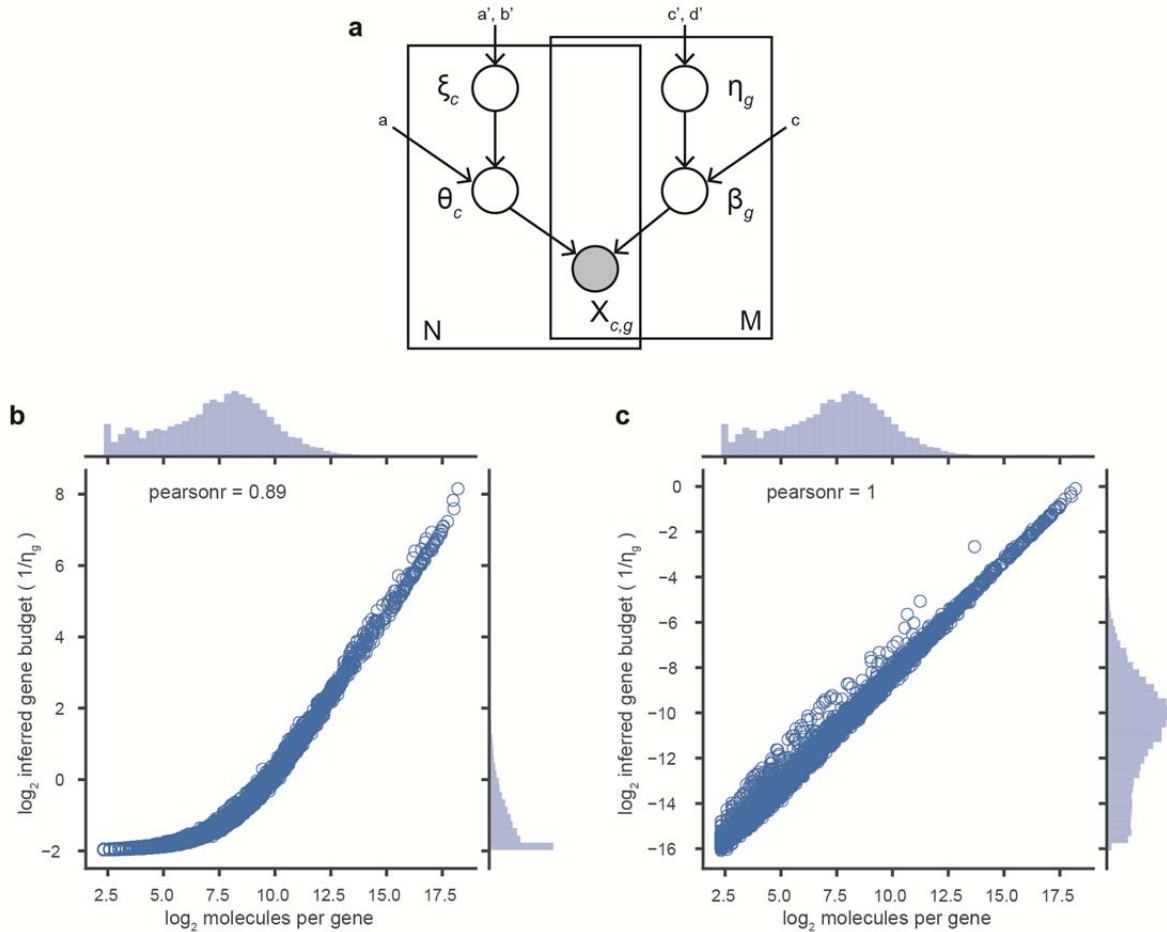


	PBMCs	Matcovitch et al.	TS543	HGG
scRNA-seq platform	10x Chromium	MARS-seq	Yuan & Sims	Yuan & Sims
Number of cells	4340	3456	9924	3109 core 3000 margin
Minimum number of cell expressing gene (for filtering)	5	5	10	10
Number of protein-coding genes after filtering	13030	8086	11807	14730
Sparsity of filtered data	90%	94%	92%	93%
K	10	10	5	14

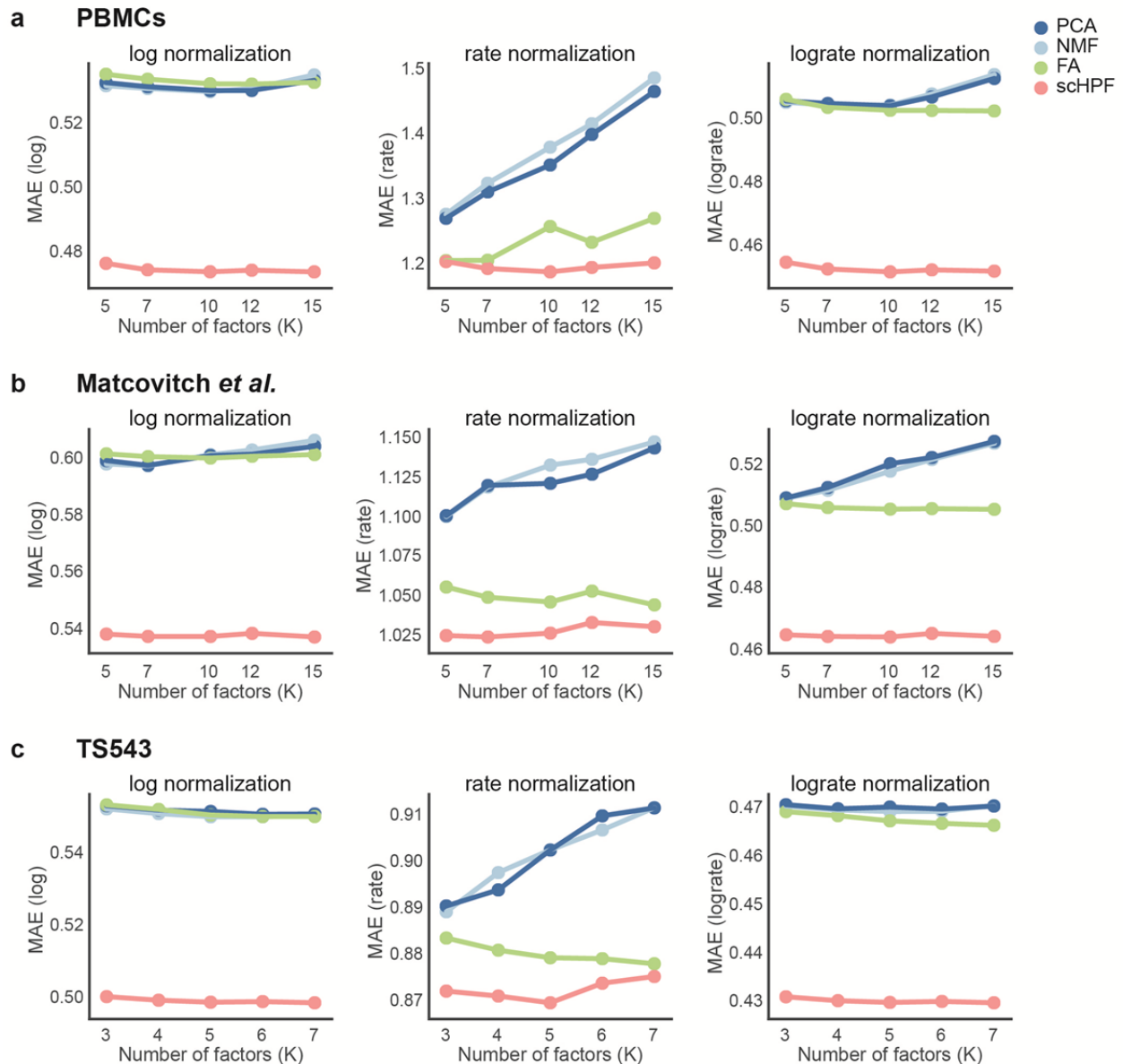
1
2
3

Supplementary Table 1: Datasets and parameters used.



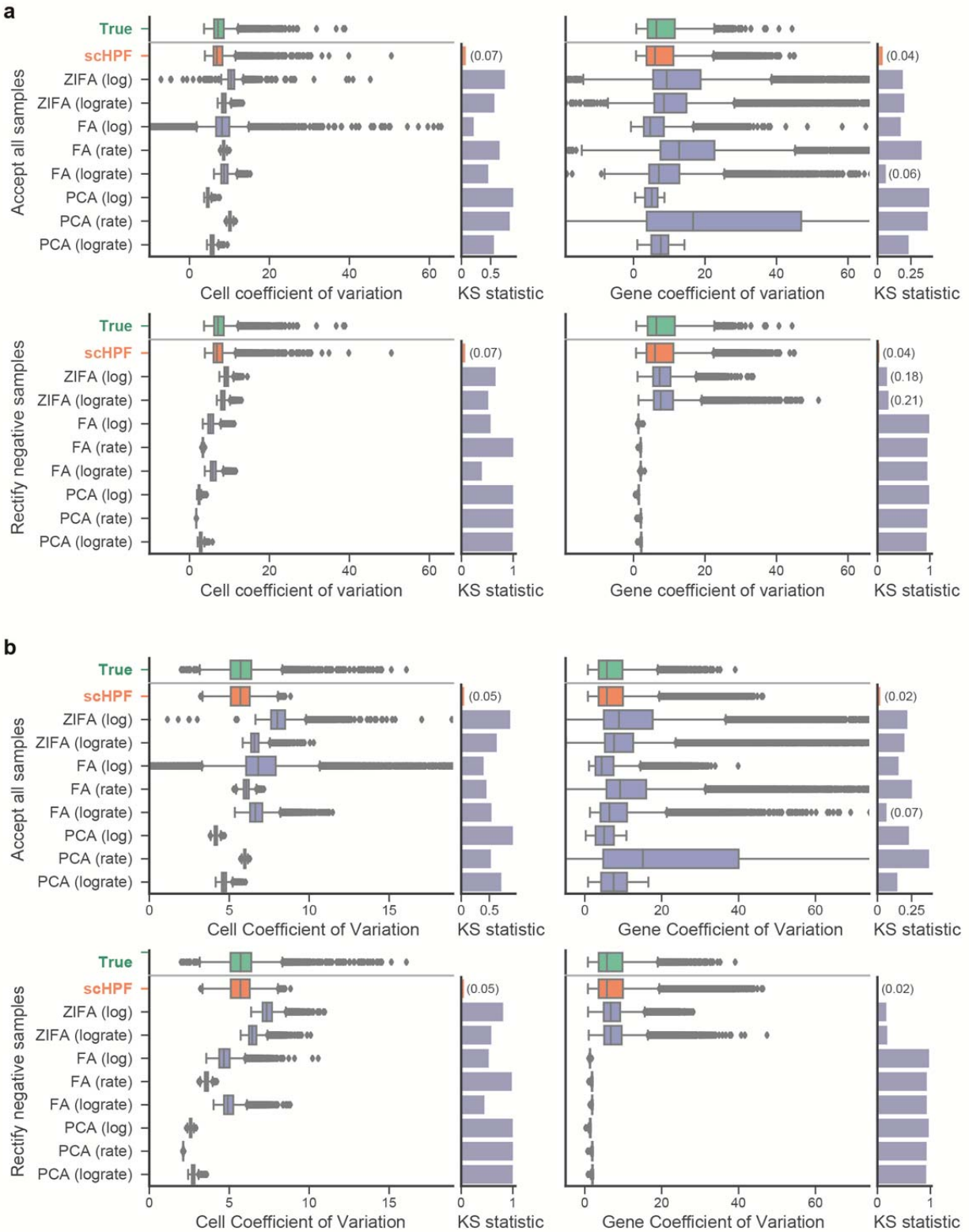
4
5
6
7
8
9
10
11
12
13
14

Supplementary Figure 1: (a) sHPF models the data matrix $X_{c,g}$ using a set of per-cell latent factors θ_c and per-gene latent factors β_g . sHPF places hierarchical priors over the latent factors through the latent variables ξ_c and η_g , which probabilistically determine the observed transcriptional output for the cell or gene. (b & c) Scatter plots of \log_2 molecules per gene (x-axes) versus the \log_2 inferred gene budgets (y-axes), with hyperparameters (b) a', b', c' and d' set to 1 or (c) determined empirically in a representative experiment on peripheral blood mononuclear cells. Histograms on top and right show the marginal probability distributions along each axis.



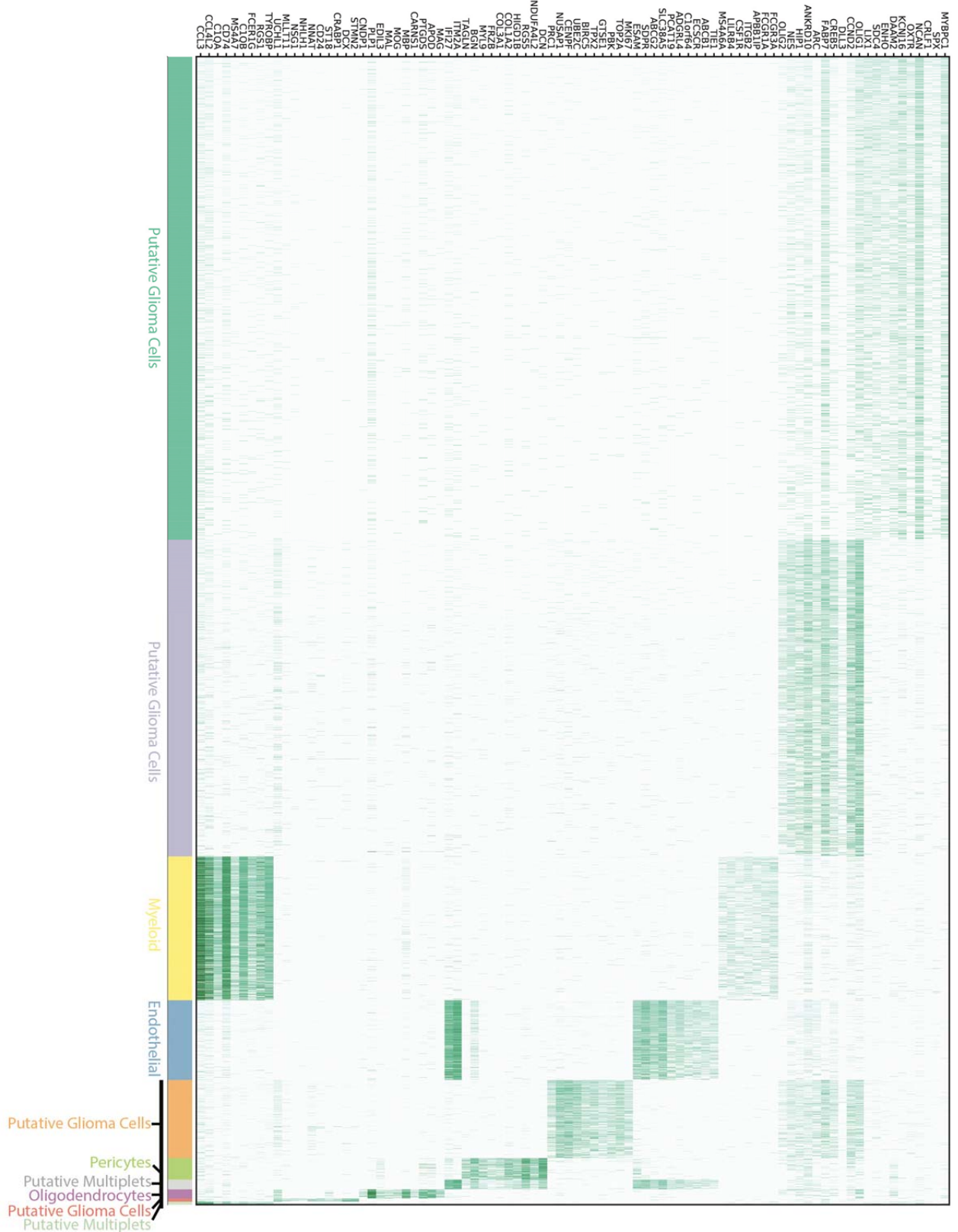
15
16
17
18
19
20

Supplementary Figure 2: Different method and normalization combinations' mean absolute error (MAE) on a withheld partition of the **(a)** PBMC, **(b)** Matcovitch *et al.*, and **(c)** TS543 datasets as compared to scHPF for several different numbers of factors. scHPF's predictions were normalized before calculating error.



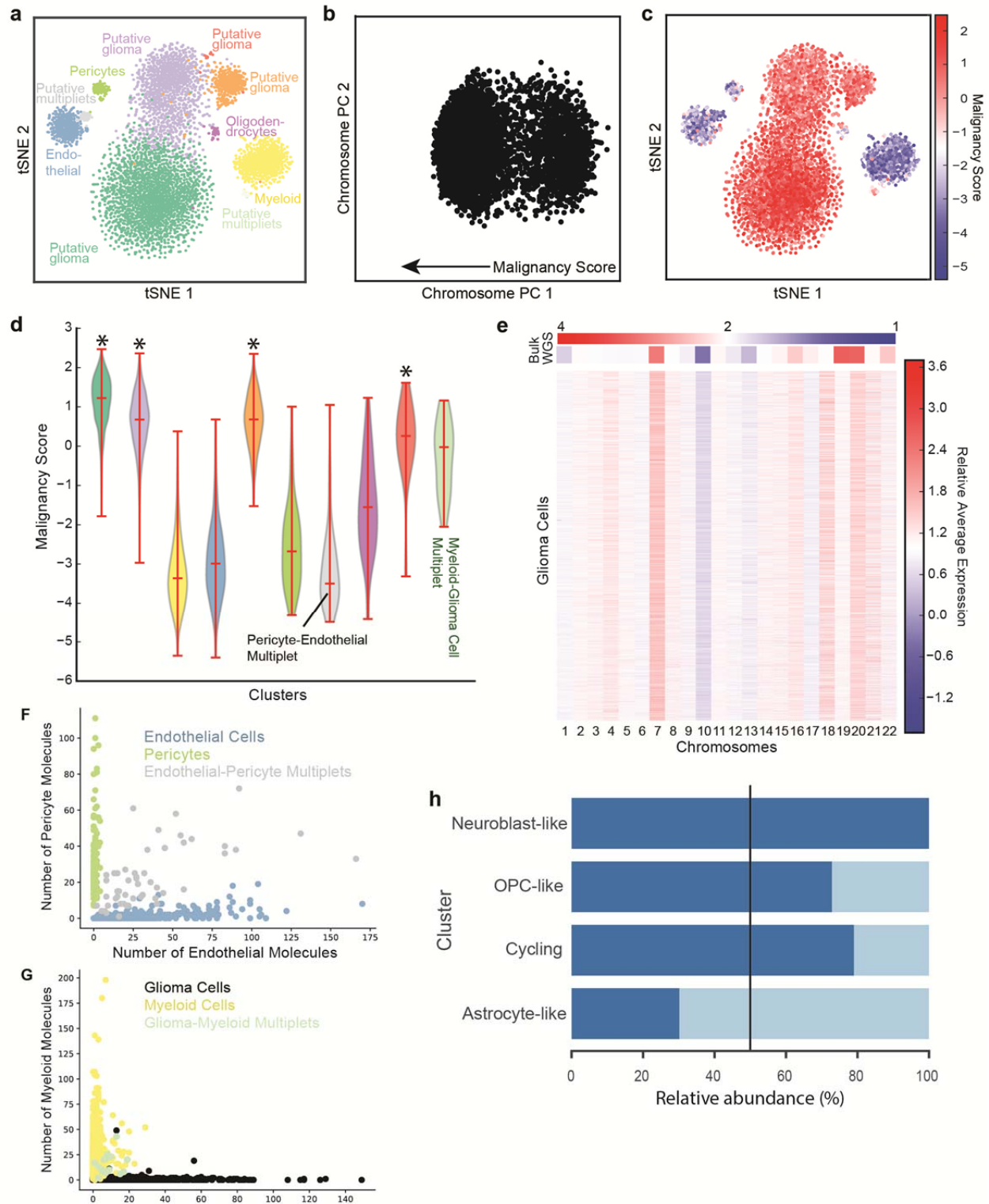
21
22
23
24
25

Supplementary Figure 3: Same as Figure 2b-c, but for **(a)** Matcovitch *et al.* and **(b)** TS543. X-axes limits for boxplots are set to include all coefficients of variation from the true distribution and scHPF, and as many coefficients of variation from other methods as possible.



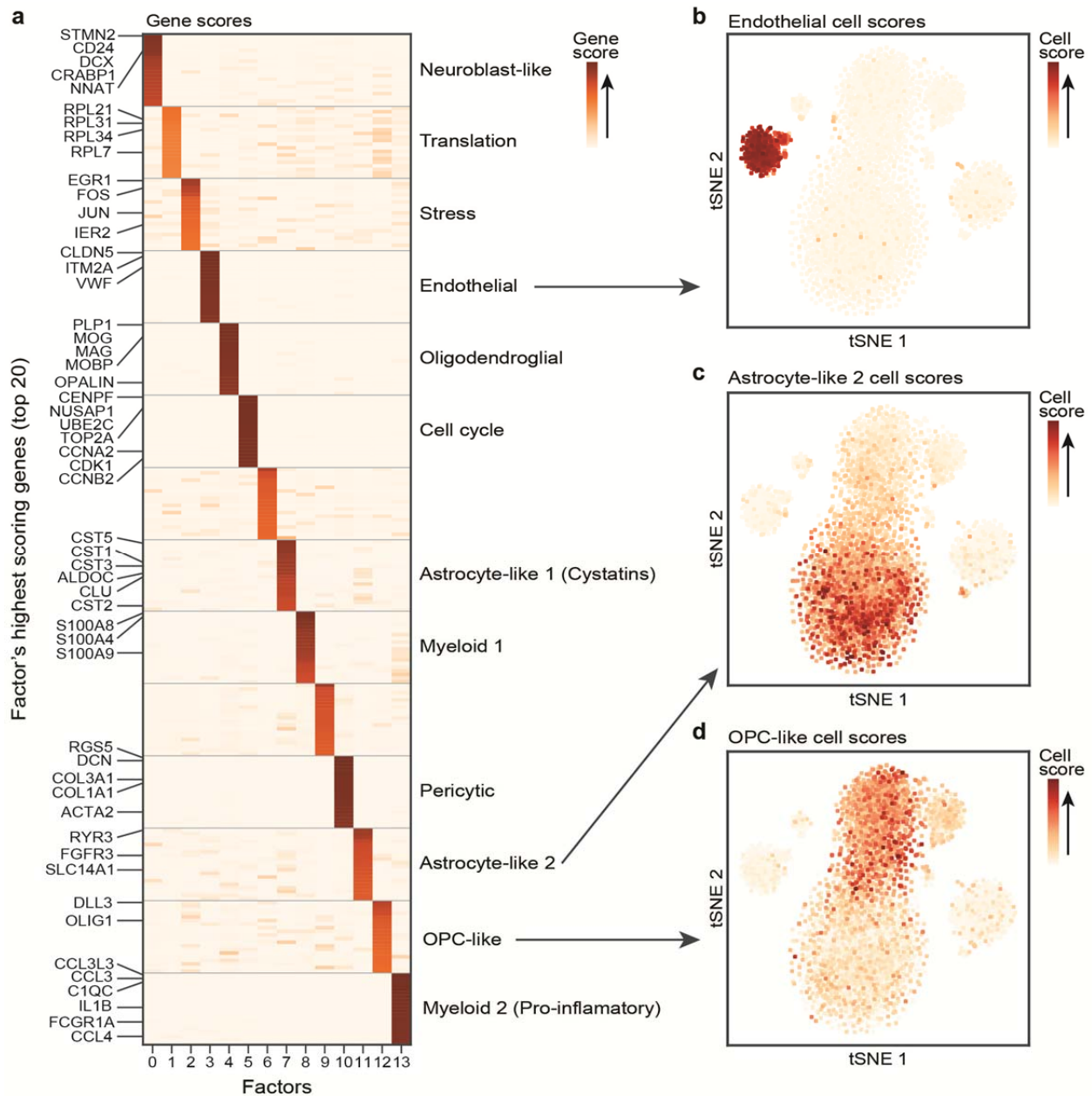
26
27
28
29
30

Supplementary Figure 4: Heatmap gene expression in a high-grade glioma with cells (columns) ordered by Louvain cluster (Methods) and genes (rows) selected as the top ten most specific genes in each cluster. Bottom color bar shows clusters and putative labels based on expression of canonical marker genes and aneuploidy analysis (see Figure S5).



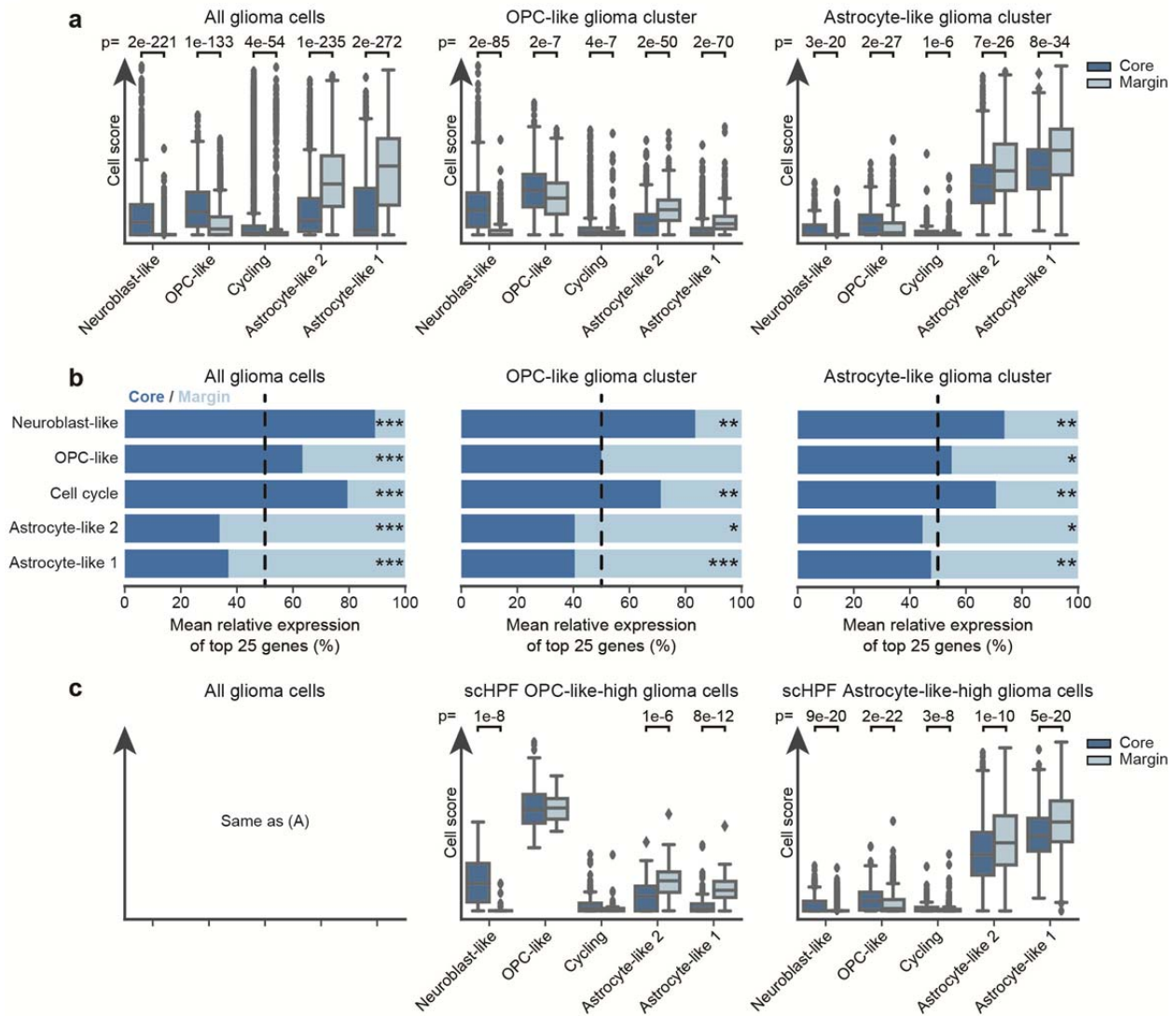
31
 32 **Supplementary Figure 5: (a)** t-distributed Stochastic Neighbor Embedding (tSNE) [1] plot of
 33 tumor cells, labeled by cluster (also see figure S4). **(b)** PCA of whole-chromosome expression
 34 for each cell. The first principle component (PC1), which we call a malignancy score, separates
 35 putative glioma from non-malignant cells. **(c)** tSNE plot of all cells, colored by malignancy
 36 score. **(d)** Violin plots of malignancy scores for each cluster. Putative glioma clusters are
 37 starred. **(e)** Main heatmap shows putative glioma cells' (rows) relative average expression of

38 each chromosome (columns). Values generally agree with bulk whole genome sequencing
39 (WGS) of the tumor (top heatmap). **(f)** Barnyard plot of cells in the endothelial (blue), pericyte
40 (green) or endothelial-pericyte multiplet (gray) clusters. Total number of molecules for the ten
41 most endothelial-specific genes by a binomial test are on the x-axis, and total number of
42 molecules for the top ten most pericyte-specific genes are on the y-axis. **(g)** Barnyard plot of all
43 putative glioma cells (black), cells in the myeloid cluster (yellow), and cells in the putative
44 myeloid-glioma multiplet cluster (green). Total number of molecules of the ten most glioma-
45 specific genes by a binomial test are on the x-axis, and total number of molecules of the ten
46 most myeloid-specific genes are on the y-axis. **(h)** Relative abundance of glioma subpopulations
47 in the core (navy) and margin (light blue).
48



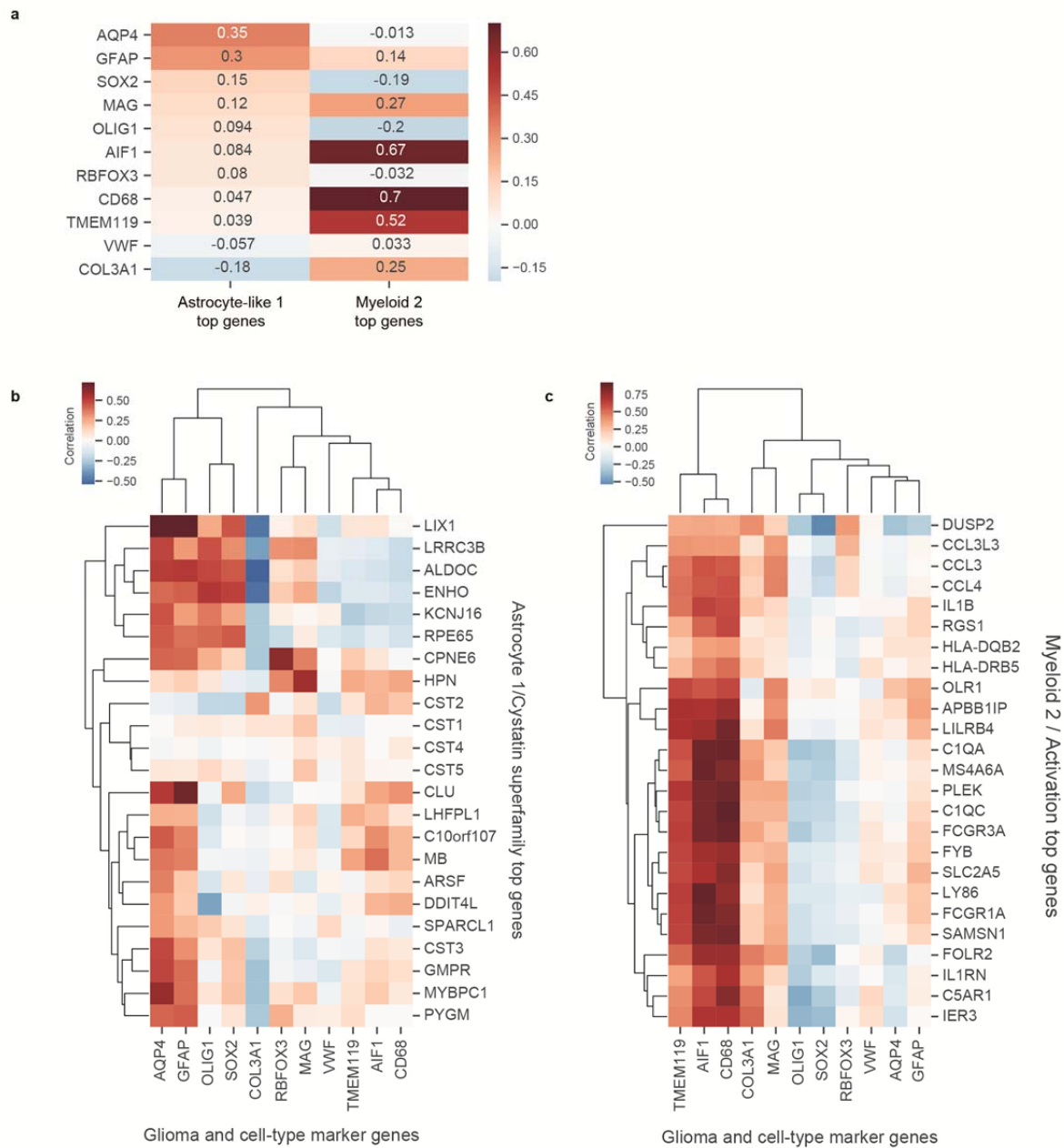
49
50
51
52
53
54
55
56

Supplementary Figure 6: (a) Heatmap of sHPF gene scores for each factor (columns) and the top twenty genes per factor (rows). Canonical marker genes and genes from a protein superfamily are highlighted. (b-d) tSNE of all cells colored by their sHPF cell scores for a factor that marks a discrete population of endothelial cells (b), one of two glioma-associated factors that highly ranks astrocyte marker genes (c), and a glioma-associated factor that highly ranks OPC maker genes.

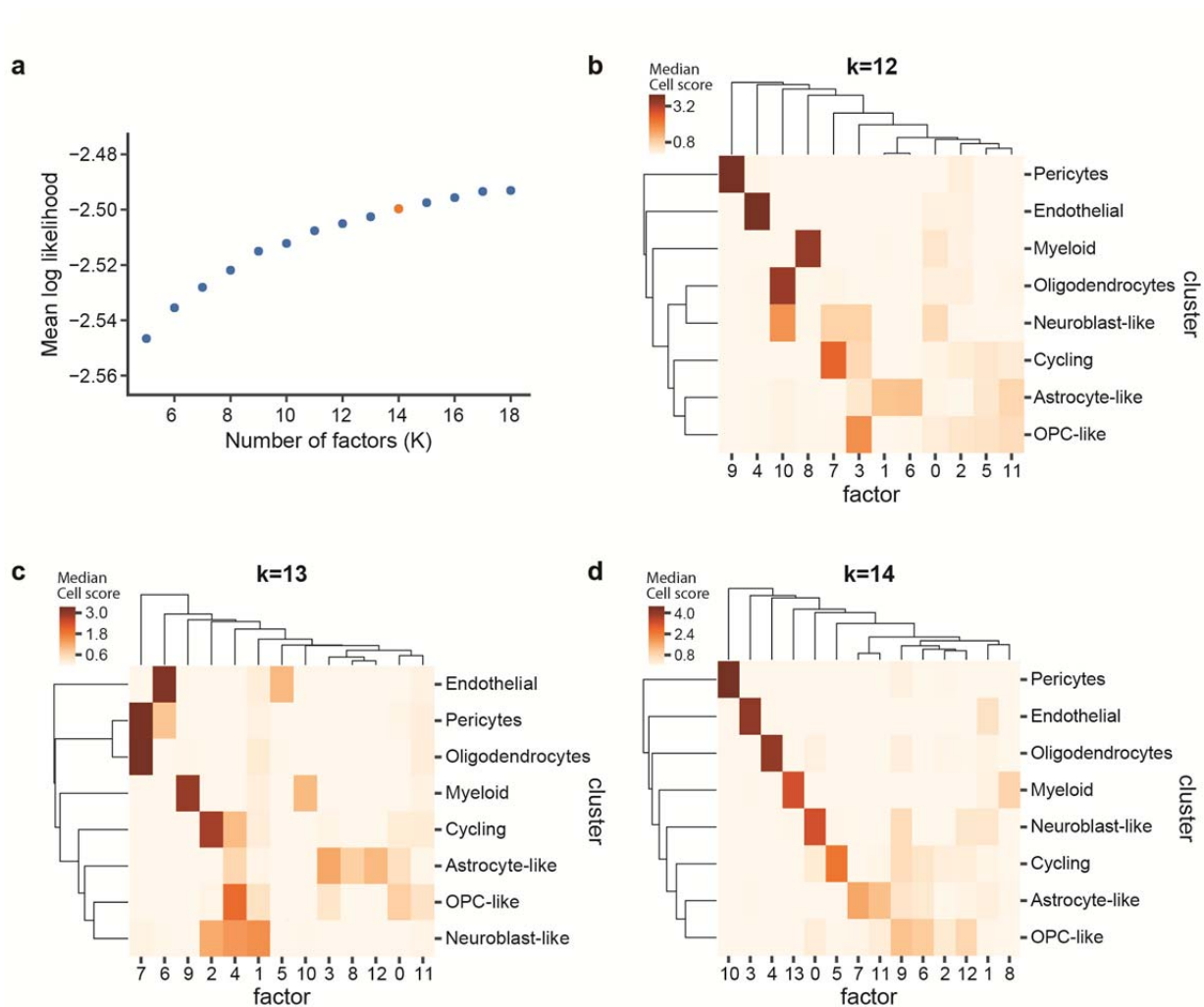


57
 58 **Supplementary Figure 7: (a)** Boxplots of scHPF cell scores for all glioma cells (left), OPC-like
 59 glioma cells (center), and astrocyte-like glioma cells (right) show strong regional bias towards
 60 the core (navy) or margin (light blue). Bracketed values show Bonferroni-corrected p-values
 61 from the Mann-Whitney U-test for the difference between two distributions. **(b)** Program scores,
 62 derived as the mean relative expression of the top 25 genes in each factor, recapitulate cell
 63 scores' regional biases. *** = $p < 10^{-50}$, ** = $p < 10^{-10}$, * = $p < 10^{-2}$. All p-values are Bonferroni
 64 corrected. Expression values were converted to counts per median and log10 scaled before
 65 averaging. **(c)** Same as (a), but with OPC-like and astrocyte-like glioma subpopulations defined
 66 as cells with maximal scHPF cell scores in the OPC-like factor or one of the two astrocyte-like
 67 factors, respectively.

68
 69
 70



71
 72 **Supplementary Figure 8: (a)** Median correlation of the top 25 genes from two scHPF factors
 73 with glioma and cell type marker genes in TCGA GBM RNA-seq. scHPF Astrocyte-like 1 is best
 74 correlated the GBM-specific marker, SOX2, and astrocyte markers. In contrast, scHPF Myeloid
 75 2 is best correlated with microglial/macrophage markers. **(b & c)** Hierarchically clustered
 76 correlation of marker genes with the top 25 genes from scHPF Astrocyte-like 1 (b) and scHPF
 77 Myeloid 2 (c).
 78



79
 80 **Supplementary Figure 9:** (a) Mean log likelihood for scHPF of a high-grade glioma at different
 81 values of K (higher is better). (b-d) Median factor score in each cluster at 12, 13, and 14 factors.
 82 With 12 factors (b), oligodendrocytes and neuroblast-like cells are both most closely associated
 83 with the same factor. Similarly, with $K=13$ (c), oligodendrocytes and pericytes are both most
 84 closely associated with the same factor. At $K=14$ (d), all clusters are most closely associated
 85 with at least one unique factor.

86
 87

88 1. Maaten, L.v.d. and G. Hinton, *Visualizing data using t-SNE*. Journal of machine learning
 89 research, 2008. **9**(Nov): p. 2579-2605.

90