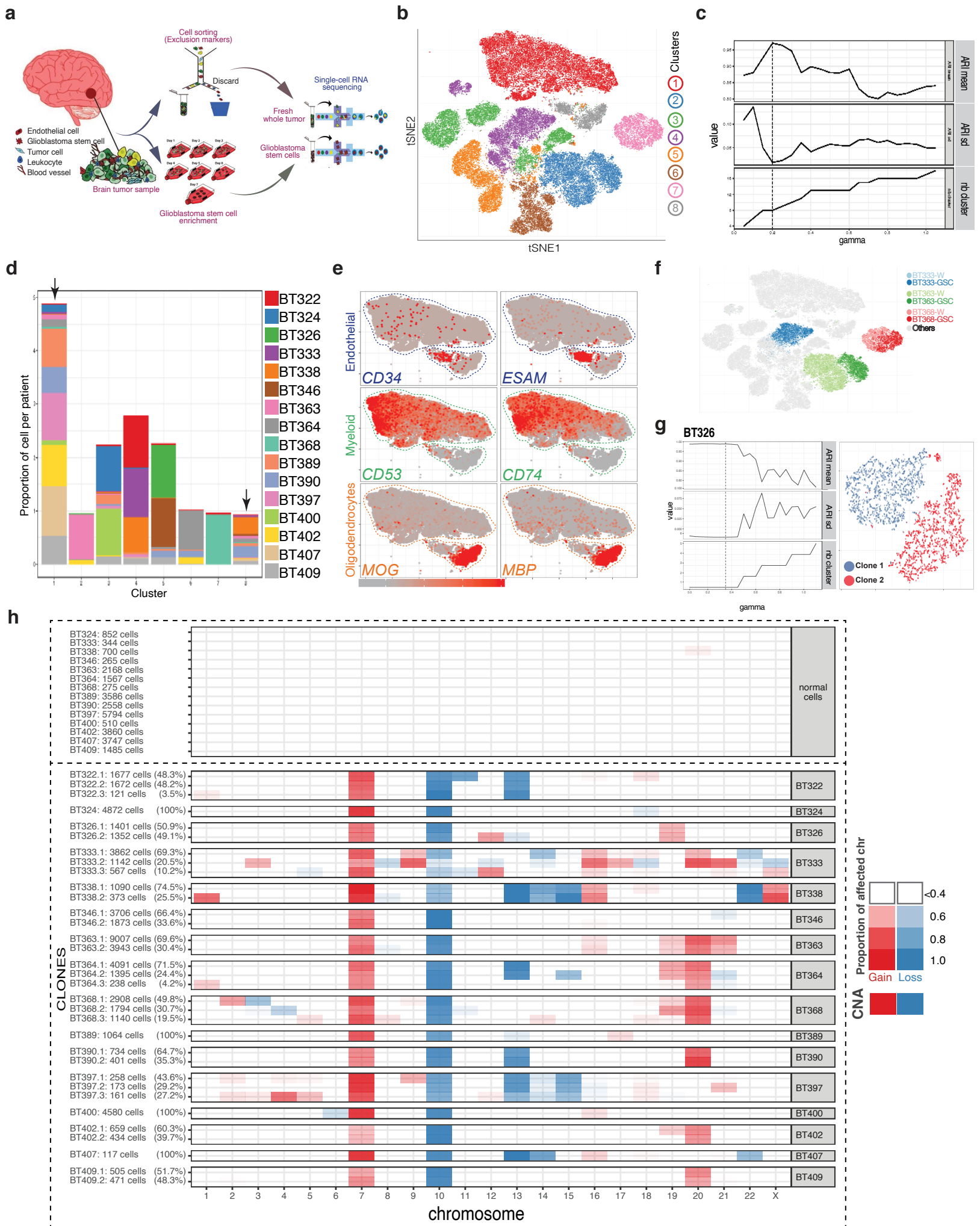


Single-cell RNA-seq reveals that glioblastoma recapitulates a normal neurodevelopmental hierarchy

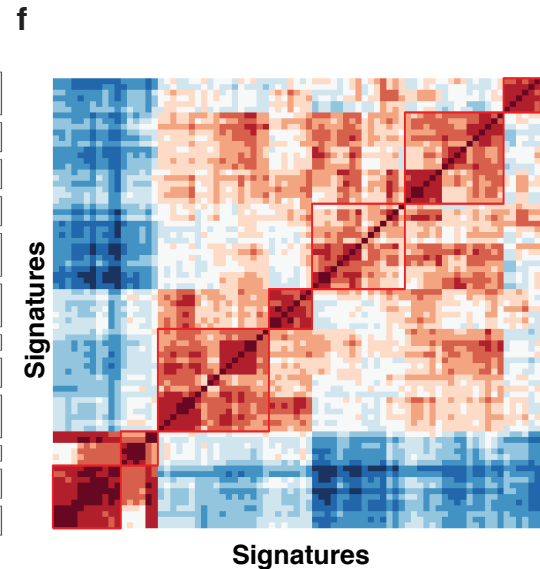
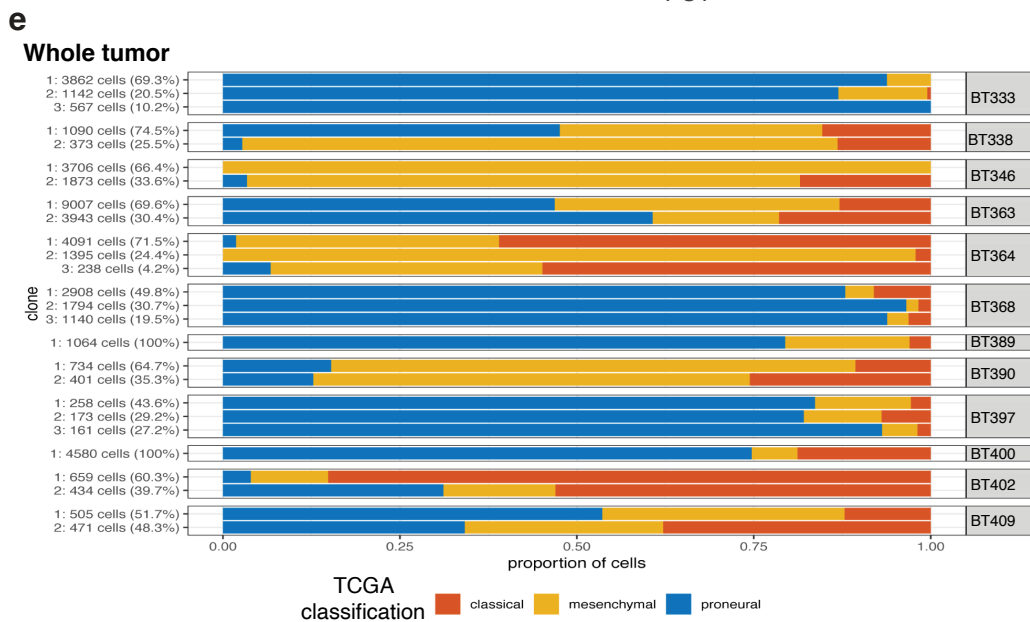
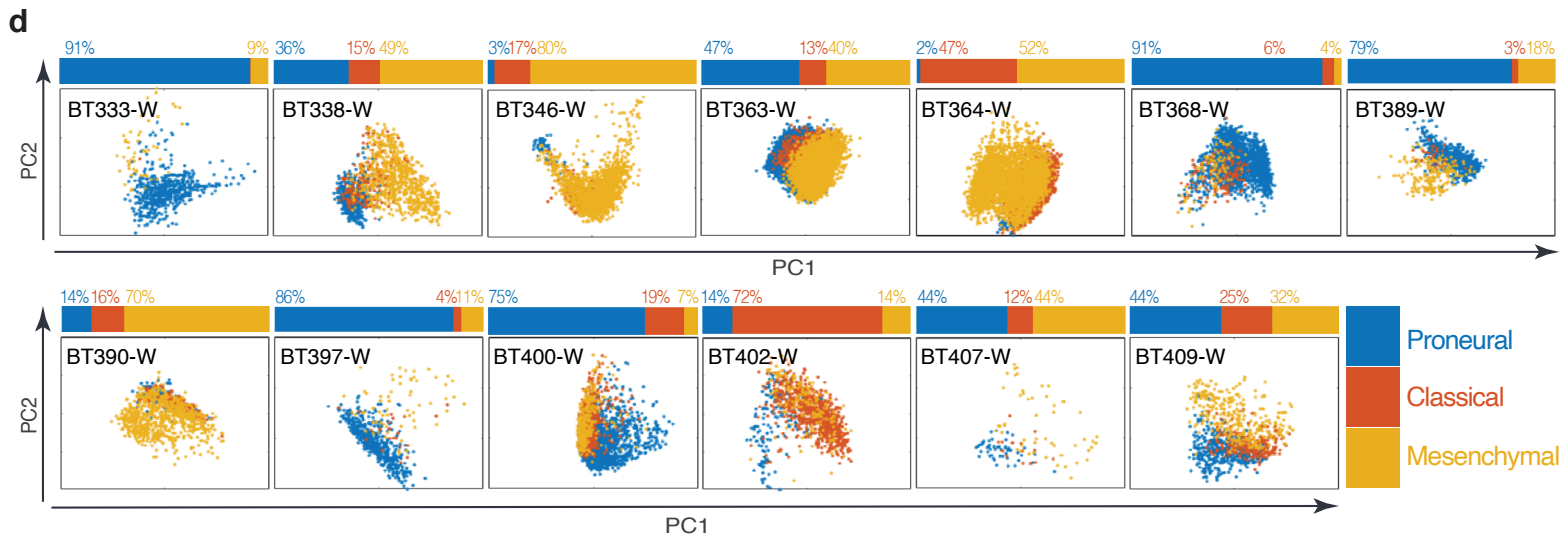
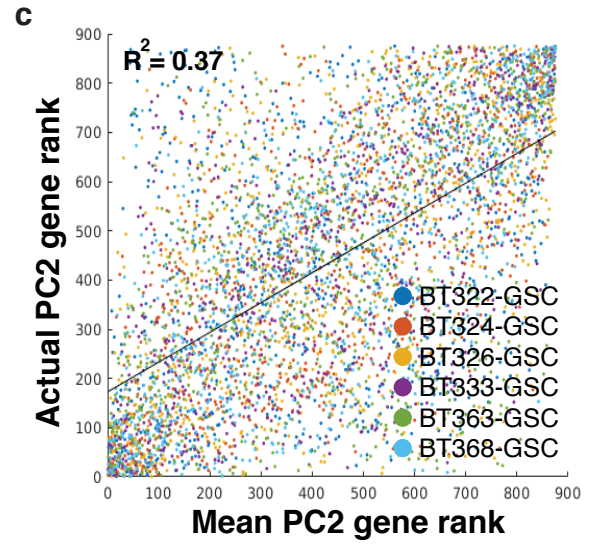
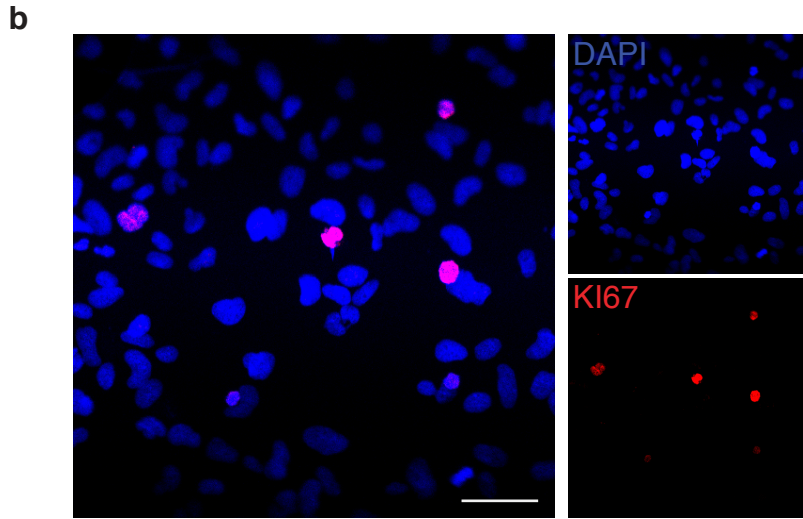
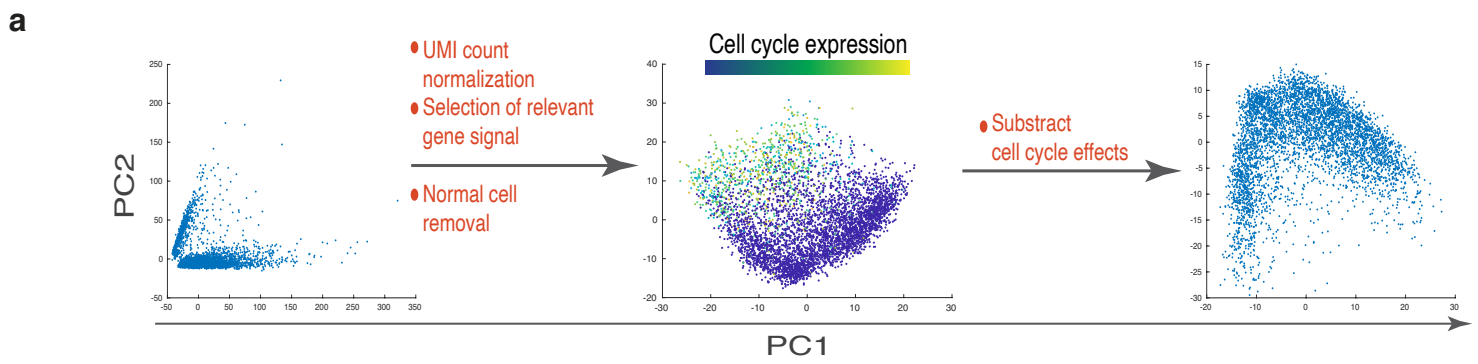
Couturier et al.

Supplementary Information



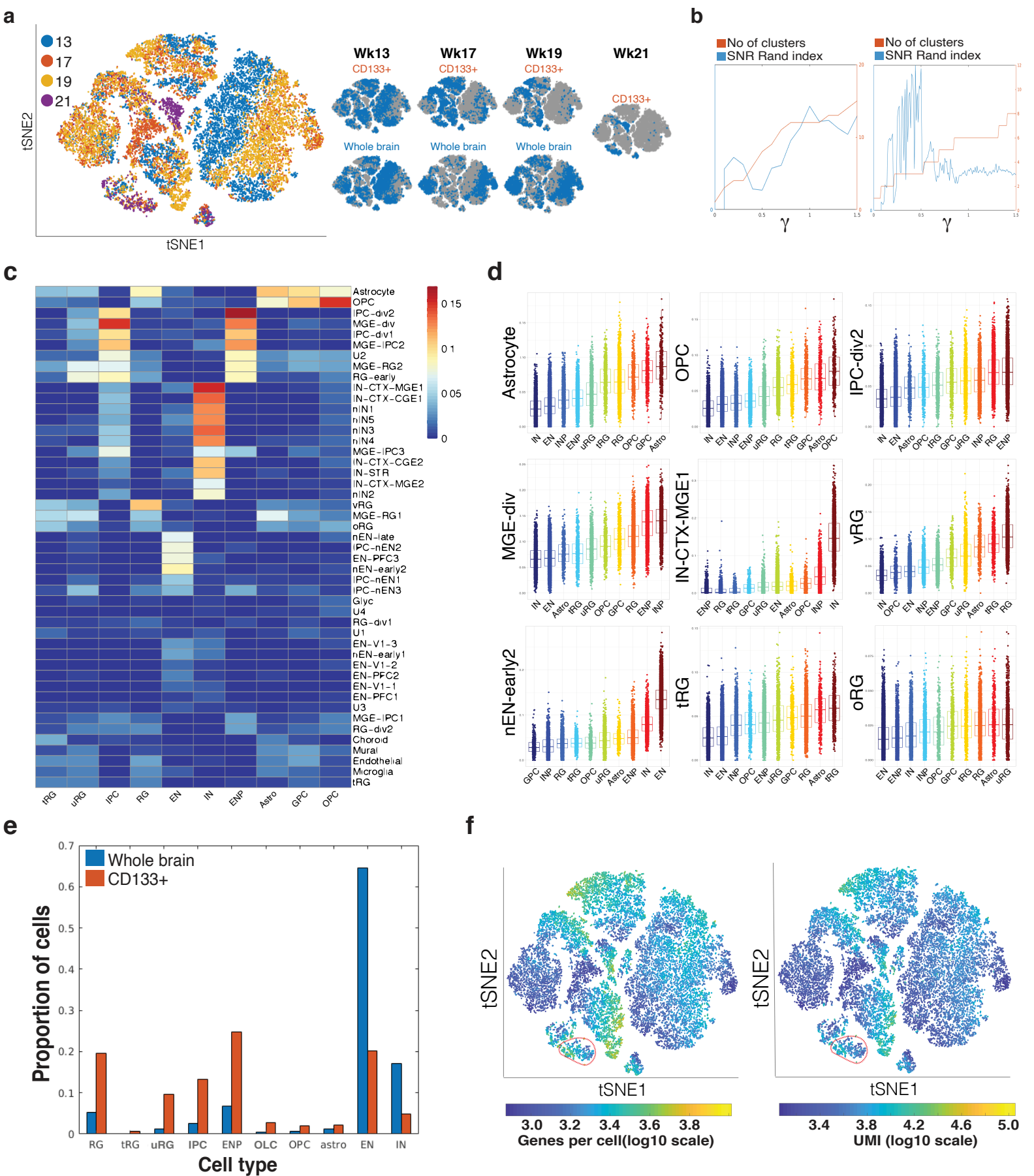
Supplementary Figure 1 Experimental design and data processing.

(a) Workflow. Brain tumour samples were extracted and cells were freshly dissociated. These cells were either used for single-cell RNA-seq (whole tumour samples) or were cultured for 7 days to produce glioma stem cell enriched samples, which were then analyzed by single-cell RNA-seq. (b) t-distributed stochastic neighbour embedding (tSNE) of location-averaged transcriptome for all tumour cells (whole tumour and glioma stem cells) from 16 patients. Cells are colored by cluster. (c) Rand index by resolution parameter (γ) for the clustering of location-averaged transcriptomic data. Top plot shows the mean Rand index, middle plot shows associated standard deviation, and bottom plot shows the median cluster number obtained for the associated γ value. (d) Proportion of patient cells for each cluster seen in Supplementary Fig. 1b. Arrows point at the clusters containing normal cells (e) Expression of endothelial (CD34 and ESAM), myeloid (CD53 and CD74) and oligodendrocyte (MOG and MBP) genes in clusters devoid of CNAs (encircled clusters in Fig. 1a). (f) tSNE of location-averaged transcriptome for patients with paired GSC and whole tumour. (g) Example of determination of clonal population in a tumour. Rand index by resolution parameter (γ) for the clustering of location-averaged transcriptomic data. Left: top plot shows the mean Rand index, middle plot shows associated standard deviation, and bottom plot shows the median cluster number obtained for the associated γ value. Right: t-distributed stochastic neighbour embedding plot for this data with the cells colored by cluster/clone. (h) Copy number aberrations (CNA) heatmap for all main clones and normal cells of each patient. The transparency shows how much of the chromosome is affected, starting at 50%. All findings shown were significant ($p < 0.001$) using the Wilcoxon test (two-sided).



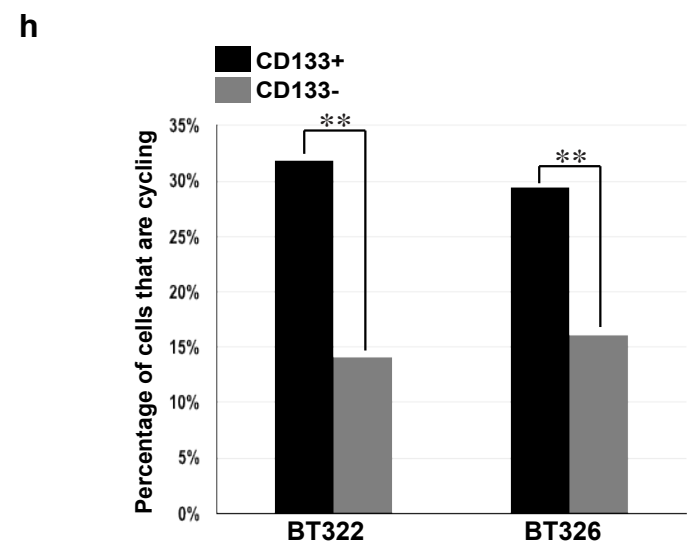
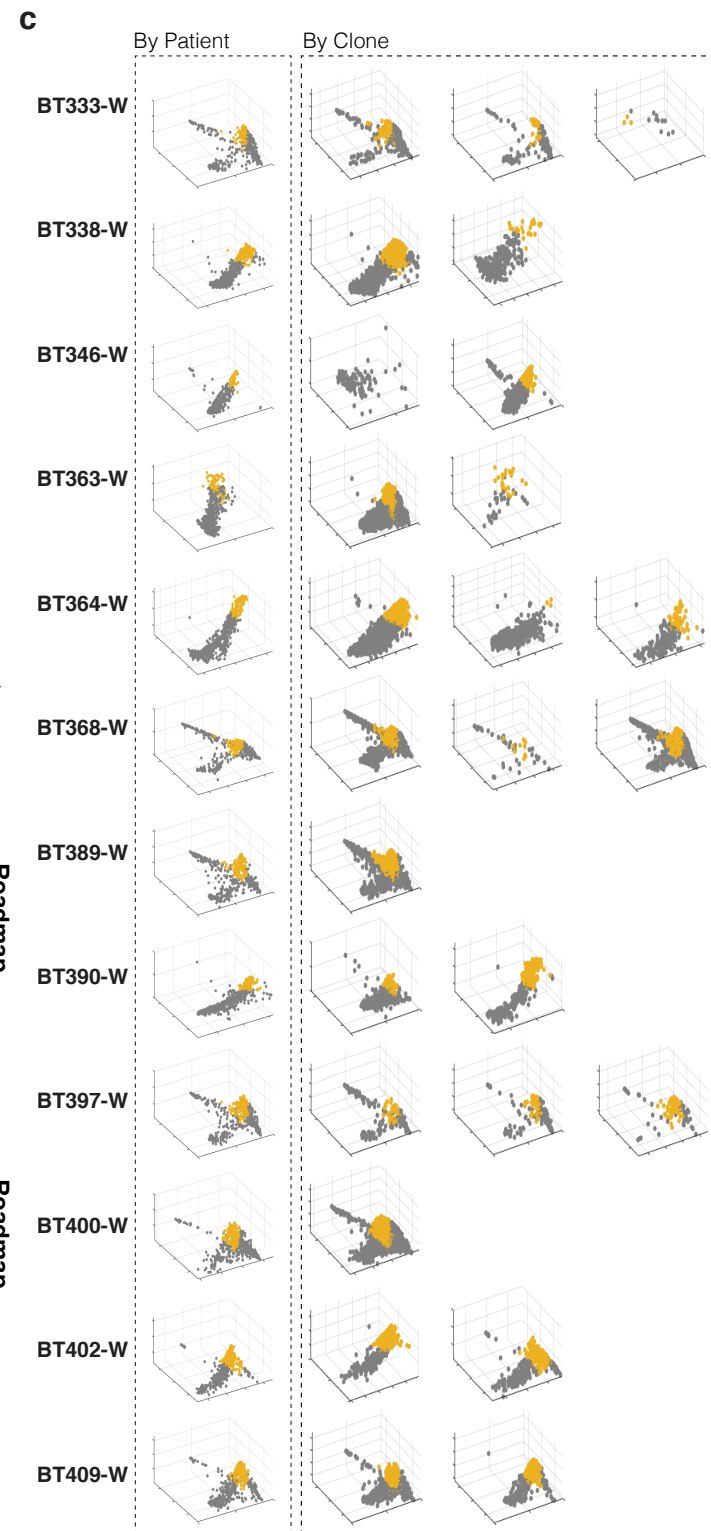
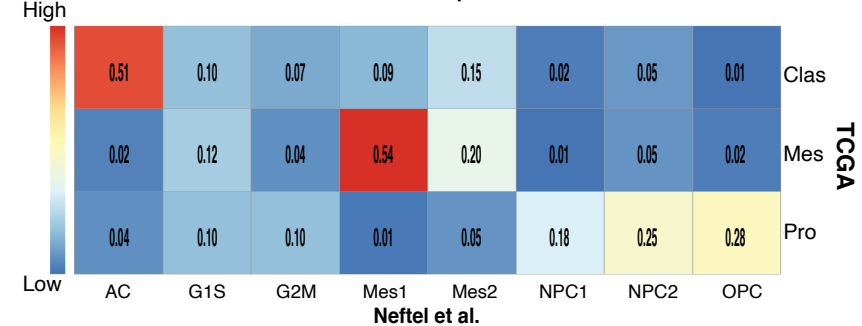
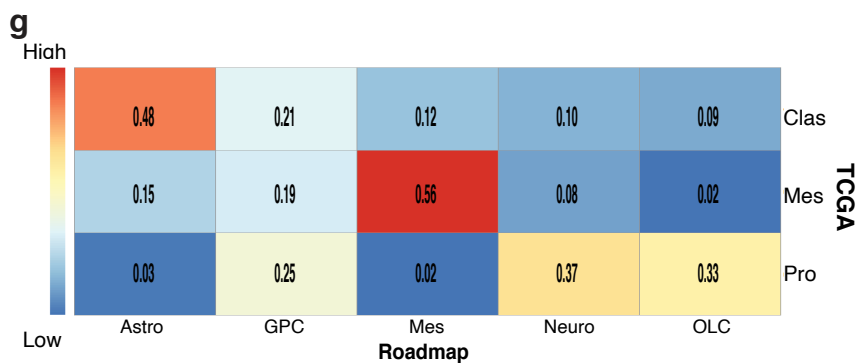
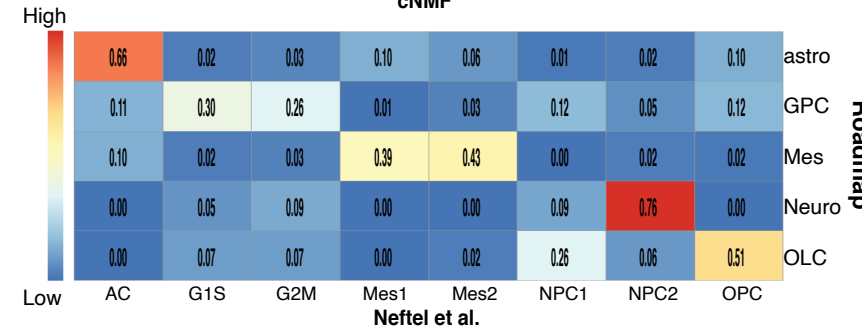
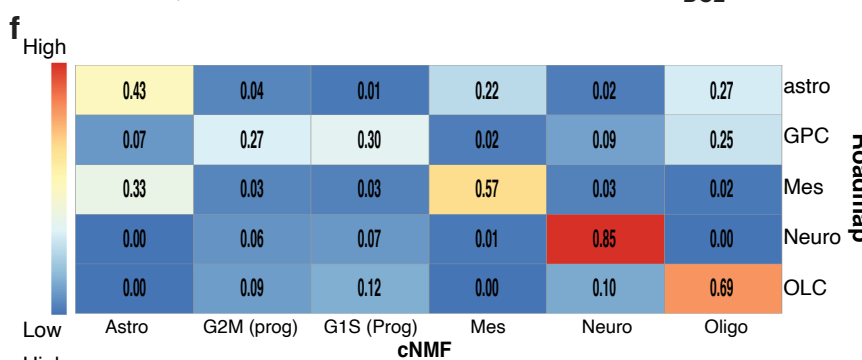
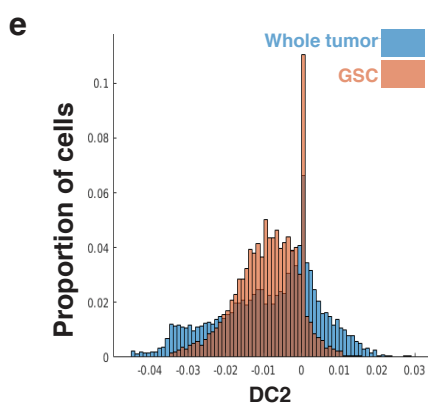
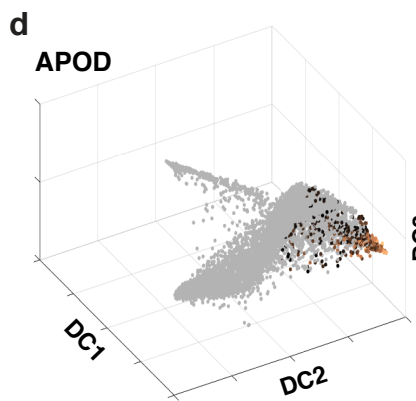
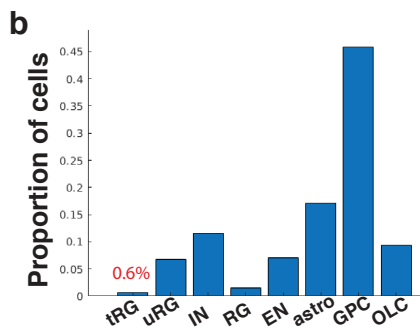
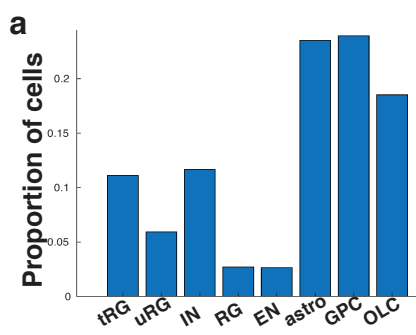
Supplementary Figure 2 Heterogeneity in glioblastoma.

(a) Bioinformatics workflow for normalization of gene expression data, gene selection, and removal of the cell cycle effect. Representative data for each step from BT324-GSC is shown by PCA. Raw data (left plot) is processed to remove low quality or normal cells and normalized for UMI count. Highly variable genes are selected. The resulting plot (middle) is characteristic for cell cycle effects as shown by the labelling of MKI67 expression in the cells. Removal of cell cycle effects yields the right plot. UMI: unique molecular identifier. (b) Representative immunofluorescence images showing Ki67 expression in enriched glioma stem cells after 1 week in culture. Image capture was done using a 20x objective with a further 2x digital magnification. n=6 biologically independent GSCs were sampled and stained. Scale bar: 50 μ m. (c) Expected and actual rank of genes by PC2 correlation. The actual gene rank (y-axis, one point per sample) correlates with the expected gene rank (x-axis) in all patients. (d) Cell-cycle corrected transcriptome principal component analysis of whole tumour cells for each patient. Cells are colored according to their TCGA subtype. Most TCGA subtypes were detected in each patient. (e) TCGA subtype by clone for each whole tumour sample. Differences in TCGA subtypes from clone to clone within a patient sample are statistically significant ($p < 0.01$, Chi-Square test) in approximately 90% of the cases. (f) Heatmap showing the correlation between signatures obtained for all 13 whole-tumour samples. Signatures are ordered by according to the hierarchical clustering in Fig. 1e.



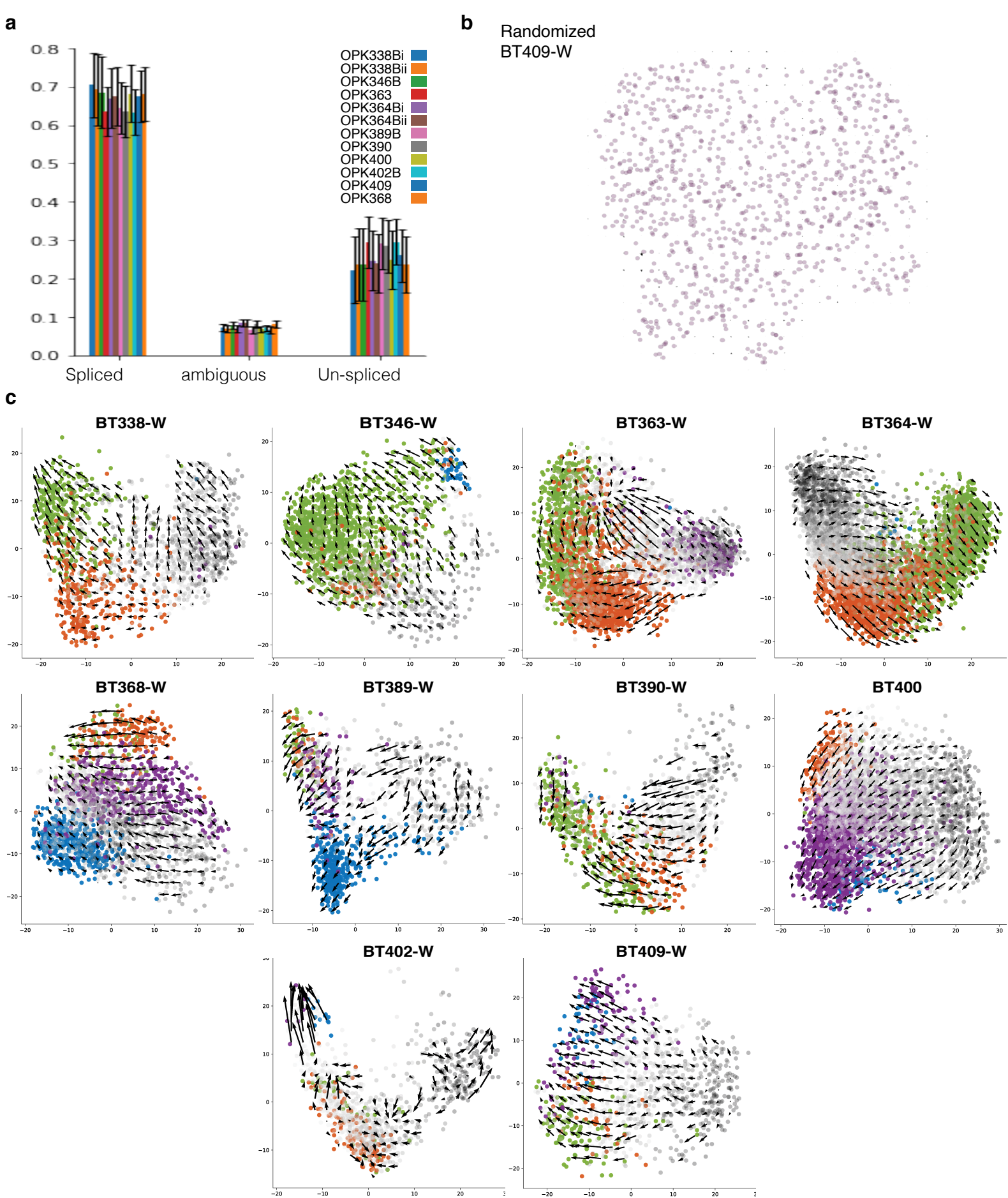
Supplementary Figure 3 Clustering by cell type in whole and CD133-positive fetal brain cells.

(a) t-distributed stochastic neighbour embedding (t-SNE) map for the fetal dataset separated by gestational age (weeks) and CD133-sorting status. (b) Rand index signal-to-noise-ratio (SNR, blue) and number of clusters (red) as a function of resolution parameter (γ) for the full fetal dataset (left) and the glial population only (right). Peak SNRs value were chosen as a final clustering solution. (c) Heatmap showing similarity (Jaccard Coefficients) between reference clusters (rows), with clusters identified in the present study (columns). (d) Box plots showing the ranked-based similarity score (see Methods) between each cluster identified in the present study (Fig. 2a) and its corresponding top 9 (out of 47) reference clusters by similarity. The box plots were ordered from lowest to highest median score (highest similarity) and represent first quartile, median, and third quartile. (e) Proportion of cells by cell type for the whole brain and CD133-positive fetal samples. The GPC subgroup is enriched following CD133-sorting. (f) t-SNE of entire fetal dataset depicting the number of genes and UMI detected per cell. The circled cell group corresponds to GPCs.



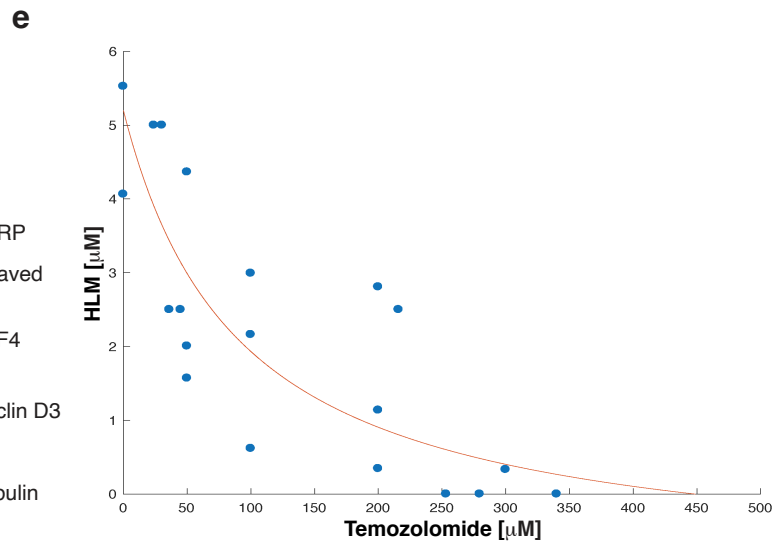
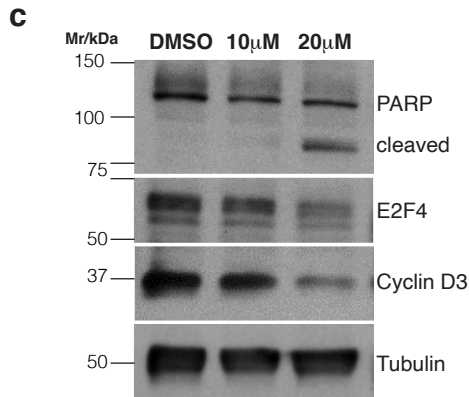
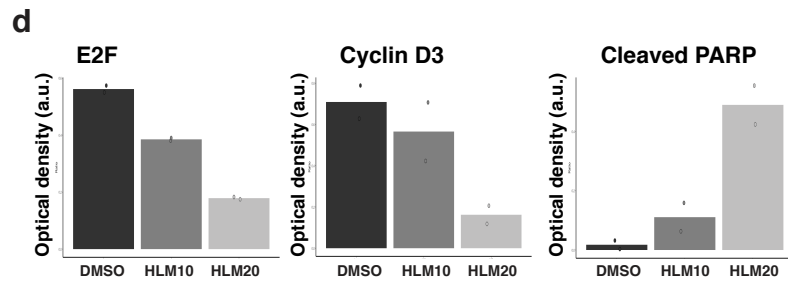
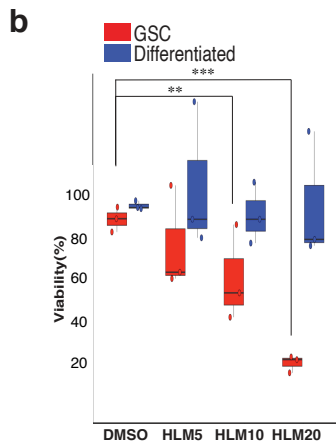
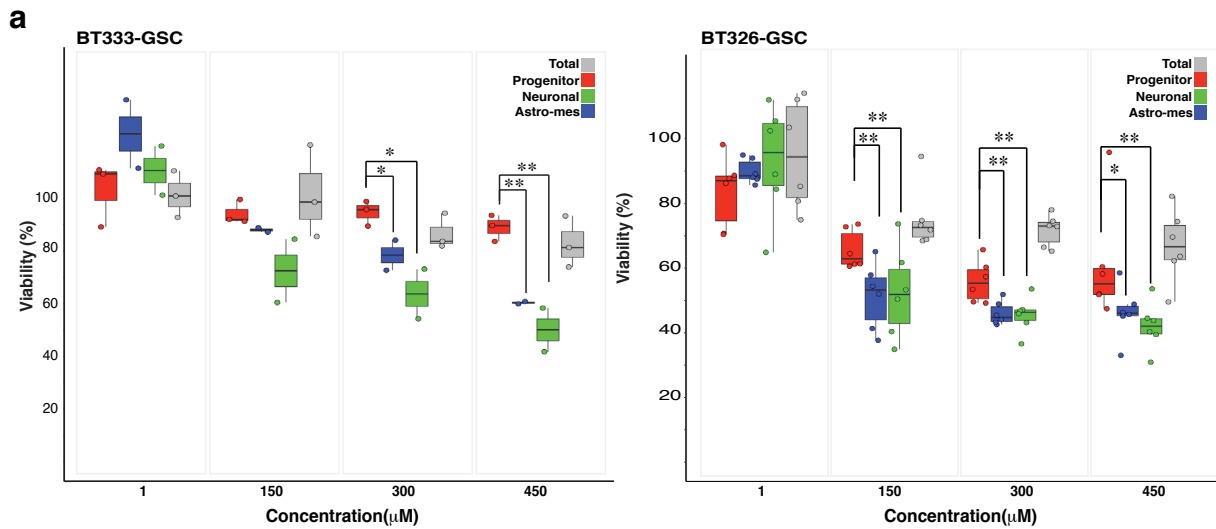
Supplementary Figure 4 Building a fetal roadmap to understand glioblastoma heterogeneity.

(a) Proportion of whole tumour cancer cells captured (see Results and Methods) by each fetal cell type. tRG: truncated radial glial; uRG: unknown radial glia; IN: interneurons (progenitors and more differentiated); RG: radial glia; EN: excitatory neurons (progenitors and more differentiated); astro: astrocytes; GPC: glial progenitor cells; OLC: oligo-lineage cells. (b) Proportion of glioma stem cells (GSCs) captured (see Results and Methods) by each fetal cell type. The tRG cell type captures a very low proportion of GSCs. See abbreviations in (a). (c) Whole tumour cancer cells isolated from 12 patients and mapped onto the roadmap in diffusion component space, by patient and by clone when applicable. Cells in yellow correspond to glial progenitor cancer cells. (d) Whole tumour cancer cells projected onto the roadmap showing expression of APOD. (e) Histogram of DC2 values for glioma stem cells (GSC) and whole tumour cells. A shift of cells from lower values (astrocytic) to intermediate values of DC2 occurs with GSC enrichment. Intermediate values of DC2 correspond to high values in DC3 and the glial progenitor cancer cell subtype. (f) Similarity analysis in our dataset of the fetal roadmap-based glioblastoma cell types, cNMF derived glioblastoma signatures (Fig. 1e), and the cell states describe by Neftel et al.⁴⁴ (g) Comparison of the TCGA subtypes to the cell type signatures in the TCGA dataset. Cell type signatures were obtained using the roadmap and Neftel et al. TCGA subtypes were obtained using Gliovis. (h) Tumour immunolabeling, using Ki67 as a marker of cell proliferation, shows that the percentage of cycling cells in the CD133-positive population is significantly higher than that of CD133-negative population in two patients. Chi-square test: p values are 0.003 (GSC322, n=272 cells from 5 fields of view) and 0.006 (GSC326, n=319 cells from 5 fields of view). ** indicates $p < 0.01$.



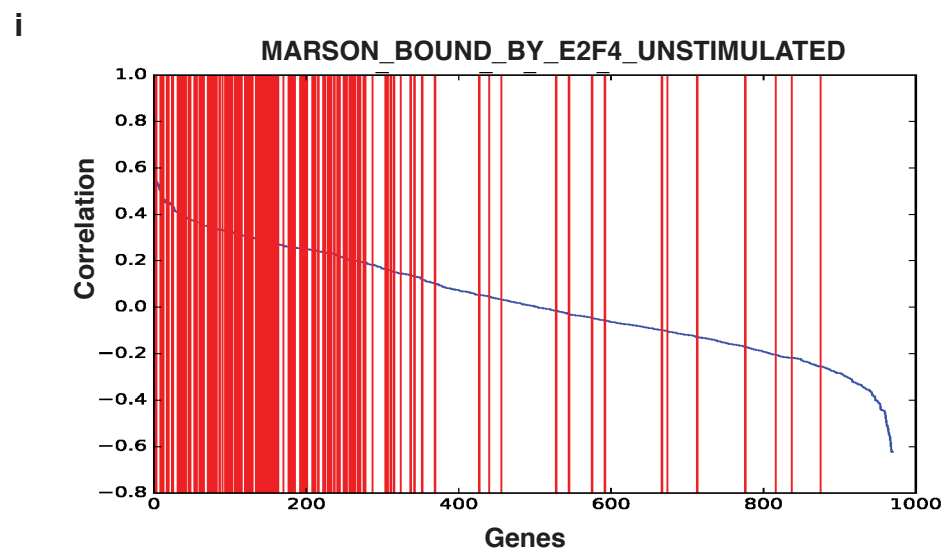
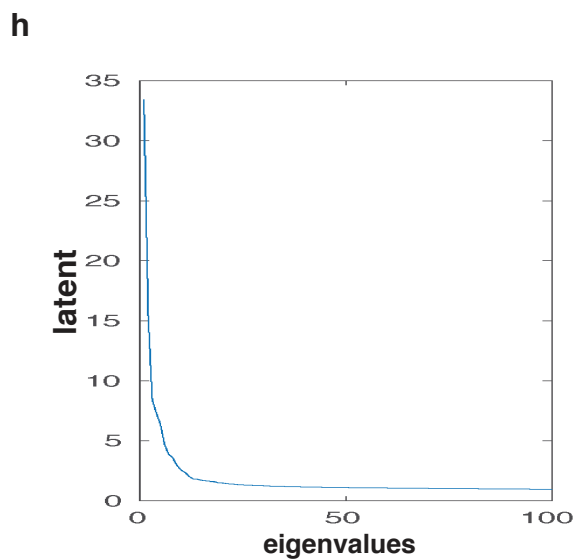
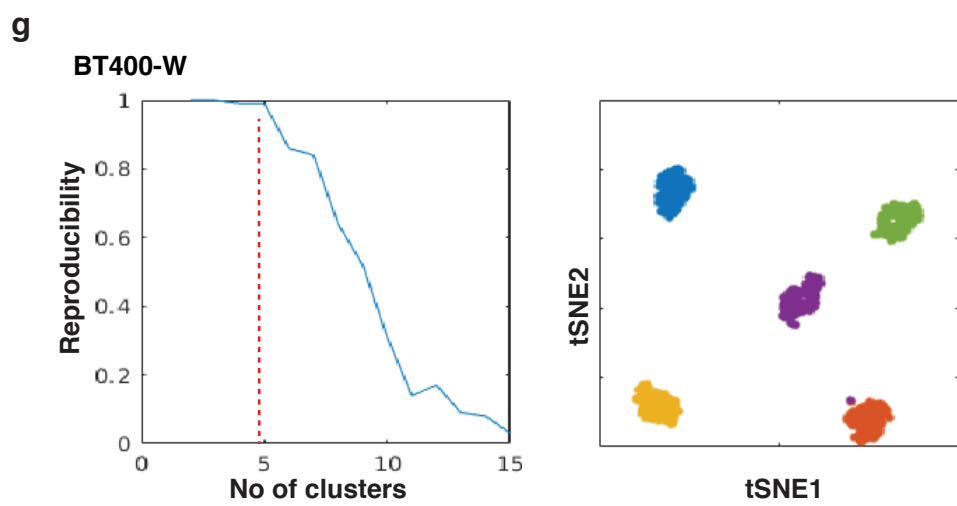
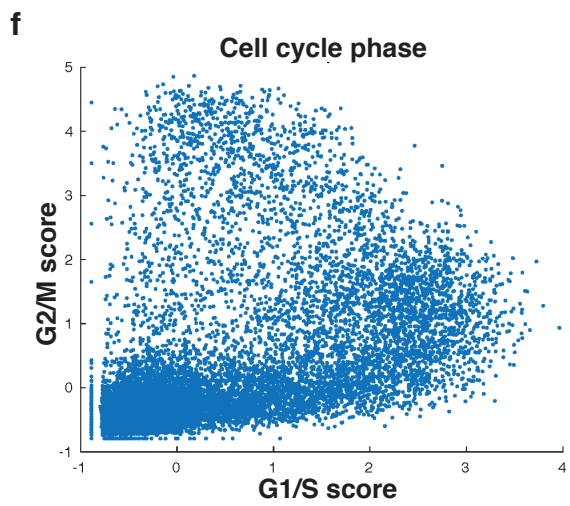
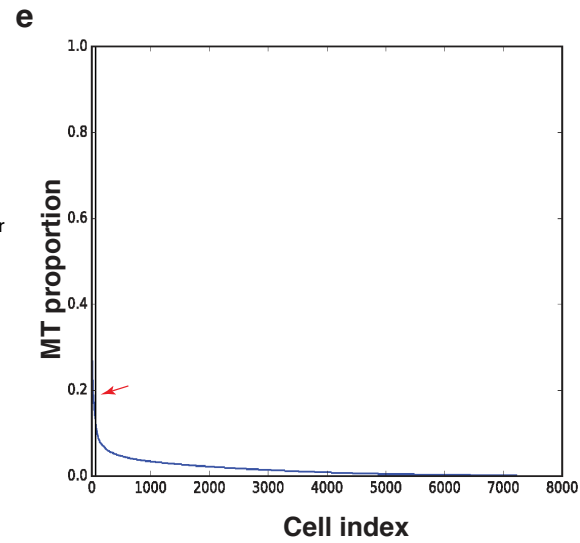
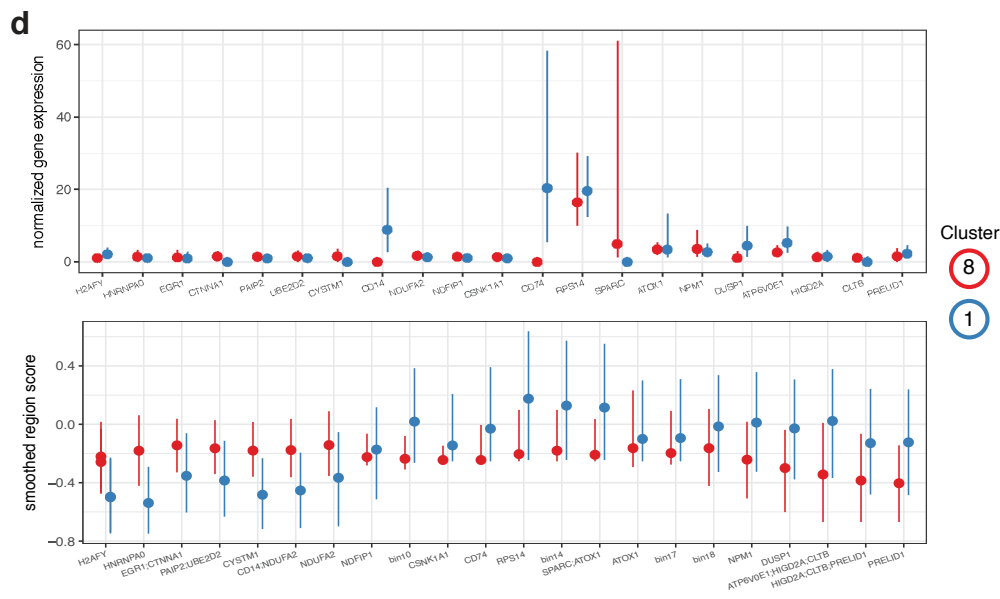
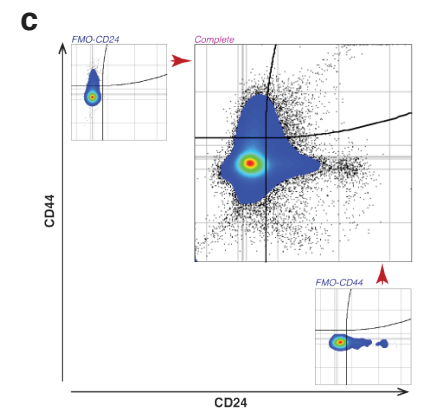
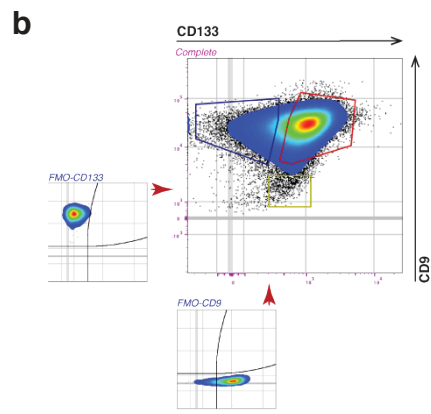
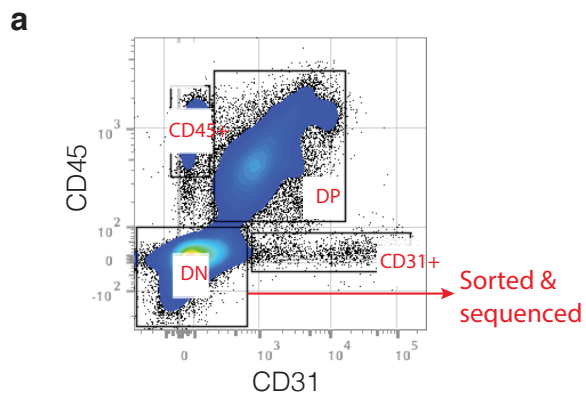
Supplementary Figure 5 RNA velocity of cancer cells.

(a) Proportions of spliced and unspliced RNA in each sample ($n=5e5$ to $5e7$ per sample depending on number of cells and reads per sample) as measured by Velocyto plotted as mean \pm standard deviation. "i/ii" indicate samples that were captured twice. All samples contain 20 to 30% unspliced molecules (b) Representative sample of a velocity map with randomized data. No clear dynamic is discerned. (c) Velocity field superimposed to PC1 and PC2 embedding of cells for all samples. Cells are coloured by cell type according to Fig 5a.



Supplementary Figure 6 Temozolomide assays show progenitor resistance and HLM006474 targets progenitors through inhibition of E2F4.

Supplementary Figure 6 Temozolomide assays show progenitor resistance and HLM006474 targets progenitors through inhibition of E2F4. (a) Box-whisker plots showing the proportion of viable glioma stem cells sorted by type followed by 5 days of temozolomide (TMZ) treatment, normalized to corresponding vehicle control. Each graph represents 1 patient sample with $n = 2-3$ biological replicates and 3 technical repeats for BT333-GSC and $n = 4-6$ biological replicates and 3 technical repeats for BT326-GSC. Biological replicates are shown as overlaid dot plot. A one-tailed, two-sample equal variance t-test was used. *** indicates $p < 0.001$, ** indicates $p < 0.01$, * indicates $p < 0.05$. (b) A representative box-whisker plot showing the proportion of viable GSCs grown in serum free conditions or serum following 7 days of HLM006474 treatment, normalized to corresponding vehicle control. Data are represented as mean \pm SE; three experimental replicates ($n = 3$) (p -values: DMSO-H10, 0.003; DMSO-H20, 4.4×10^{-12}). *** indicates $p < 0.001$, ** indicates $p < 0.01$ (one-tailed, two-sample t-test). $n = 5$ independent biological replicates were done in total with technical triplicates and showed similar trend. (c) Levels of E2F4, cyclin D3, and cleaved PARP were analyzed by Western blot from two different GSC lines; representative data are shown. $n = 2$ biologically independent GSC lines were analyzed. Full western blot found in Source Data. (d) Quantification of relative intensity of E2F4, cyclin D3, and cleaved PARP over tubulin are shown in the bar chart. $n = 2$ biological replicates for each treatment group for each experiment. (e) Isobolographic analysis of HLM006474 and temozolomide. Data show the isobole for 40% efficiency in three biological repeats. Each data point is the average of three technical repeats. Curved line indicates the null, additive isobole at 40% efficiency. P -value = 0.74 (two-tailed t-test with $n-1$ degrees of freedom). All box plots represent the first quartile, median, and third quartile with whiskers corresponding to 1.5 times the interquartile range.



Supplementary Figure 7 Cytometry and bioinformatics processing.

(a) Gating strategy for primary patient samples dissociated with anti-CD45 and anti-CD31 antibodies to remove endothelial and immune cells. (b) and (c) show the gating strategy to sort glioma stem cells (Fig. 6, Fig. 7a-c and Fig. S6a). (c) gating strategy used in Fig. 1d. CD133 and CD9 double positive cells were also sorted for CD44 and CD24 double negative cells as depicted. FMO controls were used to set the gates. (d) Top plot - gene expression by cluster for genes present over a chromosomal region. Bottom plot - location-averaged gene expression is shown for the same region. An isolated but strong expression of CD74 in cluster 1 increases the location-averaged score in that region. For each gene, the point represents the median value across all samples ($n=7$ fetal samples) in the cluster while the vertical lines represent the interquartile range (25%-75% percentiles). (e) Proportion of unique molecular identifiers detected for mitochondrial genes. The vertical line (arrow) indicates the cut-off of 0.12. (f) Cell cycle plot showing G1/S and G2/M scores in the combined fetal dataset. Most cells have a value of 0 or less in both scores. (g) Example of the determination of number of signatures by cNMF (see Methods). Left: plot of reproducibility vs number of signatures. The reproducibility drops sharply below 0.9 beyond 5 signatures. Right: t-distributed stochastic neighbor embedding plot for all 100 repetitions of 5 signatures (500 signatures in total). Five well-defined clusters are apparent, highlighting the stability of the solution. (h) PCA eigenvalues for the combined fetal dataset. The first 10 eigenvalues have the highest values. (i) Example of an E2F gene set with strong enrichment in the progenitor cancer cell population. The blue line indicates the correlation of each gene for progenitor cancer cells vs astrocytic cancer cells. Red lines show the location of genes that are included in the gene set shown. The majority of genes present in the gene set correlate positively with progenitor as opposed to astrocytic cancer cells.