

## Supplementary Information for

### High resolution mouse subventricular zone stem cell niche transcriptome reveals features of lineage, anatomy, and aging

Xuanhua P. Xie,<sup>1,2</sup> Dan R. Laks,<sup>1,2,8</sup> Daochun Sun,<sup>1,2</sup> Asaf Poran,<sup>3,9</sup> Ashley M. Laughney,<sup>3</sup> Zilai Wang,<sup>1,2</sup> Jessica Sam,<sup>4</sup> German Belenguer,<sup>5,6</sup> Isabel Fariñas,<sup>6</sup> Olivier Elemento,<sup>3</sup> Xiuping Zhou,<sup>7</sup> Luis F. Parada<sup>1,2,\*</sup>

<sup>1</sup>Cancer Biology and Genetics Program; <sup>2</sup>Brain Tumor Center  
Memorial Sloan Kettering Cancer Center, 1275 York Ave., New York, NY 10065, USA;  
<sup>3</sup>Institute for Computational Biomedicine, Department of Physiology and Biophysics; Weill  
Cornell Medicine, New York, NY 10065, USA; <sup>4</sup>Biochemistry, Cell & Molecular Biology  
Graduate Program, Weill Cornell Medicine, New York, NY 10065, USA; <sup>5</sup>Centro de  
Investigación Biomédica en Red sobre Enfermedades Neurodegenerativas (CIBERNED),  
Madrid 28031, Spain; <sup>6</sup>Departamento de Biología Celular, Biología Funcional y Antropología  
Física, Universidad de Valencia, Valencia 46010, Spain; <sup>7</sup>Institute of Nervous System  
Diseases, Xuzhou Medical University, Jiangsu 221002, PR China; <sup>8</sup>Present address: Voyager  
Therapeutics, Cambridge, MA 02139, USA. <sup>9</sup>Present address: Neon Therapeutics,  
Cambridge, MA 02139, USA.

Corresponding author: Luis F. Parada  
[paradal@mskcc.org](mailto:paradal@mskcc.org)

646 888 3781

#### This PDF file includes:

- Supplementary text
- Figures S1 to S7
- Legend for Movies S1
- Legends for Dataset Tables S1 to S3
- SI References

Other supplementary materials for this manuscript include the following:

Movies S1

Dataset Tables S1 to S3

## **Supplementary Information Text**

### **Supplemental Materials and Methods**

#### **Single cell sequencing.**

The Dropseq cell-bead collection, sample preparation, library preparation, and sequencing were performed as previously described (1). Beads were purchased from Chemgenes (#MACOSKO-2011-10) and the microfluidics chip was purchased from Nanoshift. All our reagents were purchased from the recommended sources as outlined in Macosko et al., 2015 and the McCarroll lab's online Dropseq protocol, v3.1, Dec. 2015 (<http://mccarrolllab.com/download/905/>). Our primers and oligonucleotides were purchased from IDT with the identical sequences outlined by Macosko et al., 2015. Sequencing was done on a NextSeq 500 at the Cornell Sequencing core, 64bp (R2). Alignment of reads and generation of cellular expression counts were performed with the Drop-seqAlignmentCookbookv1.2Jan2016 (<http://mccarrolllab.com/wp-content/uploads/2016/03/Drop-seqAlignmentCookbookv1.2Jan2016.pdf>) and with the Dropseq tools provided by the McCarroll lab (<http://mccarrolllab.com/download/922/>). We utilized the mm10 reference genome into which we incorporated the *CGD* transgene mRNA sequence. The transgene expression level is determined by capturing of the 3' end of the transcript that encodes the *DTR* gene. Mus\_musculus.GRCm38.82 was used as our refflat GTF. The Digital Expression NUM\_CORE\_BARCODES output was set to 5000 as we knew we had less than or equal to 2000 cells per sample.

## Single cell analysis by “Seurat”

All our analysis was performed in R version 3.3.0 with R Studio. The “Seurat” R package was used to filter, normalize, cluster and select differentially regulated genes for each cluster. Seurat package (version 1.4.0.16 and, when released, version 2.0) was utilized as described in the online clustering tutorials by the Satija lab (<http://satijalab.org/seurat>) (1). Cells were excluded if genes were detected in less than 3 cells, or had less than 200 unique genes. The resultant matrix was normalized to 10,000 transcripts, log transformed, and a regression to the number of unique molecular identifiers (UMIs) was performed before dimensional reduction. Principle components 1-19 were used to identify subpopulation clusters. We compared the lists of differential genes for each group in a cohort to select genes that were unique to a group, in other words, they were signature genes that were differentially regulated in a certain group but not in any other group. This produced two lists, differentially up-regulated genes and unique up-regulated genes. We used the differentially up-regulated gene list for gene ontology analysis by DAVID (2), (<https://david-d.ncifcrf.gov/summary.jsp>). We used the top 10 unique genes (DEGs when not enough, Table S2B & 2C) for each group to produce the heatmap with the Seurat package (Figure 3C). We performed batch correction in Seurat (version 2.0) to cluster and project both the Aged samples and the sorted two-month samples together (Fig. S7B-C). To do this we used the union of the top 2000 most variable genes for each data set and ran the canonical correlation analysis (CCA) in Seurat that identifies common sources of variation between the two datasets. For this analysis we used the first 19 dimensions and all the other settings and analyses were the same as above.

When Seurat version 3 was released, we re-did both the sorted wildtype and aged samples analyses with the latest version. As the samples were from the same tissue, preparation date, sequencing run, and there were no discernible batch effects in the sorted samples (GFP<sup>hi</sup>, GFP<sup>lo</sup>, GFP<sup>-</sup>, Unsorted) we did not use the integration method for this analysis but rather simply merged the samples into one matrix before the Seurat pipeline so as not to introduce any noise from unnecessary processing. In contrast, for the aged samples we used the anchor integration method (Stuart et al., 2019). In both instances, the original results were conserved and proved robust across different pipelines and versions of Seurat. One exception was the co-clustering of H1 with H2 into one cluster in the Seurat version 3 processing of the sorted two-month SVZ samples. As the results and clustering were largely conserved across different versions of Seurat, we present the original analysis with confidence that they are robust results irrespective of the Seurat versions and processing employed.

### **Pseudotime Analysis**

Pseudotime analysis was performed using the Destiny package as detailed in (3) and in the vignette at <https://www.helmholtz-muenchen.de/icb/research/groups/quantitative-single-cell-dynamics/software/destiny/index.html>. We used H0-L2 groups from Fig. S3A for our Pseudotime analysis in Fig. 4F. In Destiny, we plotted the result with our established color codes. The animation was compiled in Sequimago and labels were added in Wondershare Filmora (Version 7.8.9). Pseudotime analysis was also performed in a similar manner on the H and L groups of the aged samples that were included with the aforementioned pseudotime matrix of H and L groups in Fig. 4F (Fig. 7B).

### **Allen Brain Atlas SVZ localization of unique genes**

The gene list used for this experiment is the same as those used in the aforementioned Heatmap section (see above). We collected images from the Allen Brain Atlas (4), <http://mouse.brain-map.org/>) where the SVZ was at its largest in the sagittal sections. We authored a macro in ImageJ that got rid of the white space in our quantification. We then made 5 sections of the SVZ from Ventral to Dorsal and used ImageJ to quantify the in situ hybridization of RNA for each of the five sections. We then normalized the values to the most ventral section and plotted the values, performed linear regression, and in GraphPad (Prism7) we calculated whether the slopes were significantly different than 0. The p-values for this test are reported.

### **Venn Diagrams**

Venn diagrams were created by InteractiVenn (<http://www.interactivenn.net/>) (5) and reproduced in Photoshop CS6 (Version 13.0.6) to improve the quality and size of the associated text.

### **Transcription factor analysis**

Transcription factors related in a sign sensitive manner to our inputted lists of up- and down-regulated genes were identified by TFactS by target gene signatures curated from microarray gene expression data (<http://www.tfacts.org>) (5, 6).

## **String network**

Gene networks were generated by String database (<https://string-db.org>) (7) with medium confidence (0.400) to assess known and predicted interactions (physical and functional) between genes in our inputted gene sets. The PPI enrichment score p-value is a measure of the significance of enrichment for network interactions within a gene set.

## **qRT-PCR with CD95+ sorted samples**

We sorted 1000 GFP<sup>hi</sup>;CD95+, GFP<sup>hi</sup>;CD95-, bulk GFP<sup>hi</sup>, GFP<sup>lo</sup>, GFP-, or DAPI-viable cells into tubes containing 3000 Dropseq beads (Chemgenes) in 100ul lysis buffer with DTT (#15508013, Thermo Scientific) (1). The tubes were rotated at 4 degrees for 10 minutes. Then we washed the samples twice with 6xSSC (1), followed with another wash with reverse transcription buffer (Thermo Scientific). Reverse Transcription mix with Template Switch Oligo were added and incubated for 30 minutes with rotation at room temperature and an additional 90 minutes with shaking at 700 rpm at 42 degrees Celsius. We then processed the samples as Dropseq samples outlined in Macosko et al., 2015 with Exonuclease treatment, PCR for 30-35 cycles, and Agencourt Ampure XP (#A63881, Beckman-Coulter) bead purification of the resultant DNA. The purified DNA content was measured on a Qubit with high sensitivity reagents (#Q32854, Thermo Scientific). qRT-PCR reactions were performed with SYBR Select Master Mix (#4472903, Applied Biosystems) in an Applied Biosystems QuantStudio Flex6, real time PCR system. The following primers were selected from PrimerBank (8),

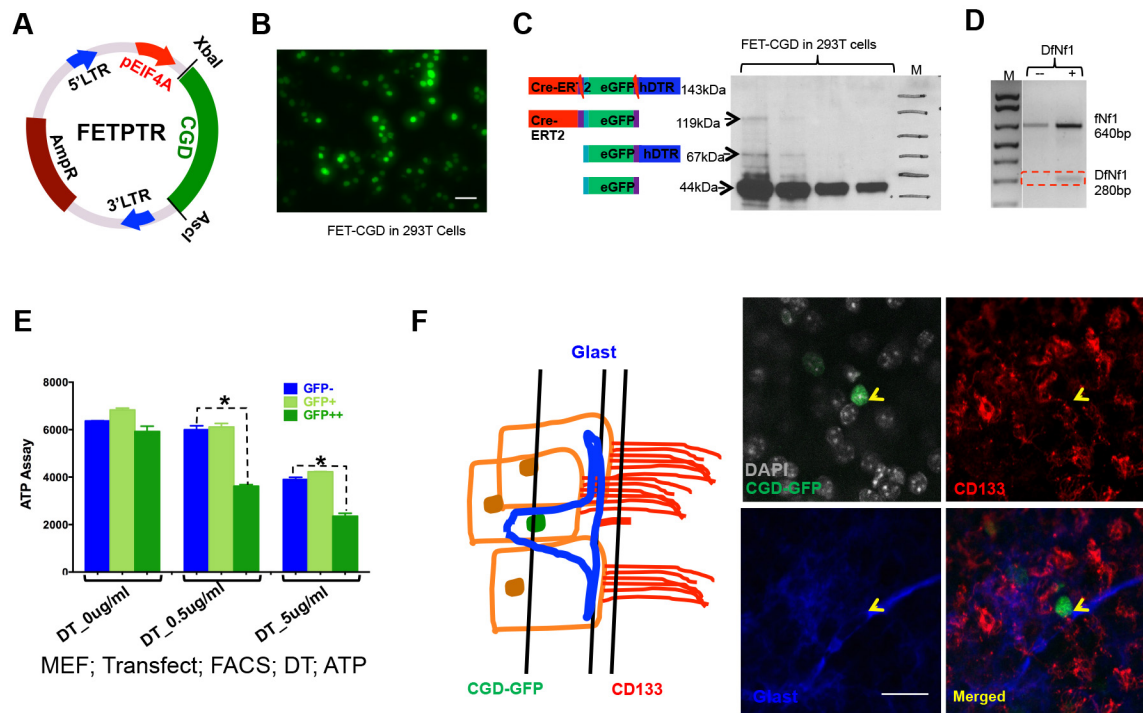
<https://pga.mgh.harvard.edu/primerbank/index.html>) except for Cd95 which was selected from Sabbagha et al. (9), and all primers were purchased from Eton Biosciences Inc.

qRT-PCR Primers for Mouse Genes

Gene (Mouse)	Fwd	Rev	Amplicon Size (bp)
1110017D15RIK	TGTCTCGGAACCGACTCT	GGAGCAAGGGCTTTCATATCC	101
Fam183b	CGTGTGGGGCAGATGAAGAAT	GGTGAATGAGGTTTCAGGAACCTG	148
Dynlrb2	GAATCCAGAGTCACAAAGGGG	GACCCGCATACTGAACCGTT	105
Rplp1	CTCGCTTGCATCTACTCCGC	AGAAAGGTTTCGACGCTGACAC	109
Tmem212	GGTACACAGGATGGAGCGTTT	GCTTCCCACAAGTGTCTCTGG	121
Actb	GGCTGTATTCCCCTCCATCG	CCAGTTGGTAACAATGCCATGT	154
Hsp90ab1	TCAAACAAGGAGATTTTCCCTCCG	GCTGTCCAACCTTAGAAGGGTC	102
Cd95 (Fas)	TATCAAGGAGGCCCATTTTGC	TGTTTCCACTTCTAAACCATGCT	148

### SCDE analysis

We employed the package SCDE (<https://www.nature.com/articles/nmeth.2967>) (10) to determine differentially expressed genes between 2 Months and 12 Months aged samples in both the combined H1-H2 groups (N=93 for 2 months, N=115 for 12 months) and in the combined H3-L0-L1 groups (N=84 for 2 months, N=63 for 12 months). We used non-log transformed, 'Seurat' normalized data matrices as the input. We chose genes with FDR P values <0.05 for further analysis including gene ontology (David, <https://david-d.ncifcrf.gov/summary.jsp>) (2).



Adapted from Mirzadeh et al., 2008

Fig. S1. This figure accompanies Figure 1. *CGD* transgene construct characterization. (A) FETPTR expression vector. (B) eGFP reporter detection in FETPTR-*CGD* construct transfected 293T cells. Scale bar: 20 $\mu$ m. (C) Western blot with GFP antibody verifies efficient production of GFP monomers. 293T cells were transiently transfected with virus packaging a pEIF4A promoter driven *CGD* cassette and collected three days later in 2X laemmli for western blot. (D) CreER protein from the *CGD* transgene can remove the floxed region within an NF1 gene allele, detected by PCR. (E) Mouse embryo fibroblasts transfected with the FETPTR-*CGD* derived viruses are sensitive to diphtheria toxin in an ATP assay. (F) Left: Diagram of classic SVZ “pinwheel” structure in which a eGFP+/Glast+/CD133+ stem cell is surrounded by CD133+ ependymal cells. The three black lines indicate the relative positions for the microscope images on the right. Right: Corresponding IHC SVZ images with a *CGD*-GFP<sup>hi</sup> nucleus and Glast+ cytoplasm surrounded by CD133+;GFP- ependymal cells. Scale bar: 20 $\mu$ m.



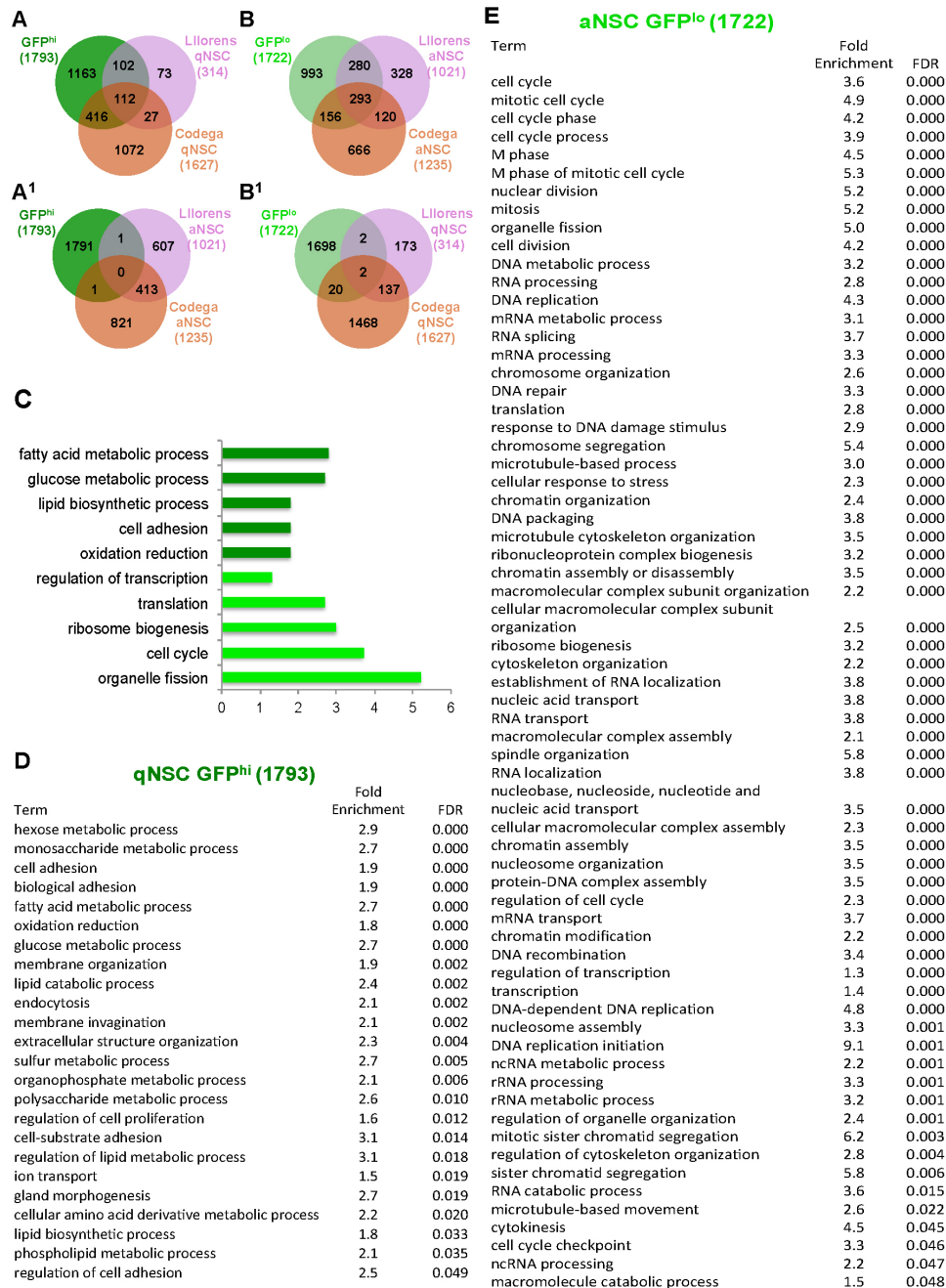


Fig. S2. This figure accompanies Figure 2. Gene expression analysis of *CGD*-GFP<sup>hi</sup> and GFP<sup>lo</sup> cells point to a quiescent versus a proliferative state respectively. (A) Venn diagram analysis comparison of *CGD*-GFP<sup>hi</sup> DEGs and published quiescent NSC signatures demonstrates significant overlap but not with activated NSCs (A<sup>1</sup>). (B) GFP<sup>lo</sup> cell signatures preferentially overlap with that of published activated NSCs but not with quiescent NSCs (B<sup>1</sup>). (C) Gene Ontology analysis of *CGD*-GFP<sup>hi</sup> and GFP<sup>lo</sup> profiles is consistent with a stem versus progenitor state. (D) and (E) Known and unknown biological processes related to quiescent or proliferative NSCs are associated with GFP<sup>hi</sup> (D) or GFP<sup>lo</sup> cells (E).

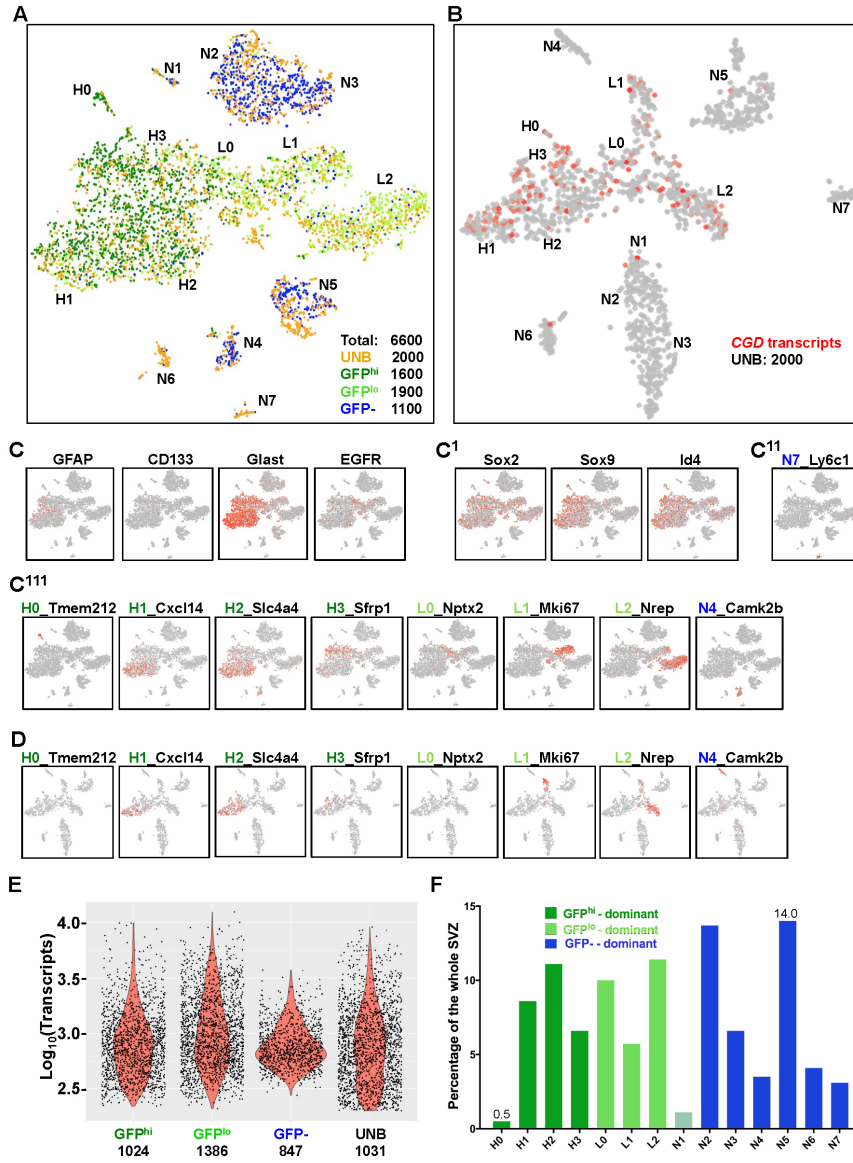


Fig. S3 This figure accompanies Figure 3 and Table 2. Drop-seq analysis of whole SVZ. (A) Combined 6600 cells yield fourteen groups that integrate the GFP<sup>hi</sup>, GFP<sup>lo</sup>, and GFP<sup>-</sup> subgroups. Note UNB cells exist in each of the fourteen groups. (B) tSNE projection of UNB cells shows similar groups and enriched CGD expression in H1-L2 lineage. (C) Distribution of sorting markers for stem/progenitor cells used in previous studies. (C<sup>1</sup>) Representative distribution of NSC and progenitor specific transcription factors; (C<sup>11</sup>-C<sup>111</sup>) candidate markers identified in this study for GFP<sup>-</sup>:N7 and different SVZ stem/progenitor groups. (D) Representative markers for NSC/Progenitor/neuron lineage facilitate the identification of respective populations in UNB cells. (E) Transcript numbers for each of the four sorted and Dropseq analyzed samples present a “bell” shape distribution. Average transcripts per cell for each sample are indicated at the bottom. (F) SVZ UNB cells contain fourteen distinct subgroups ranging from 0.5% to 14% of the whole tissue.

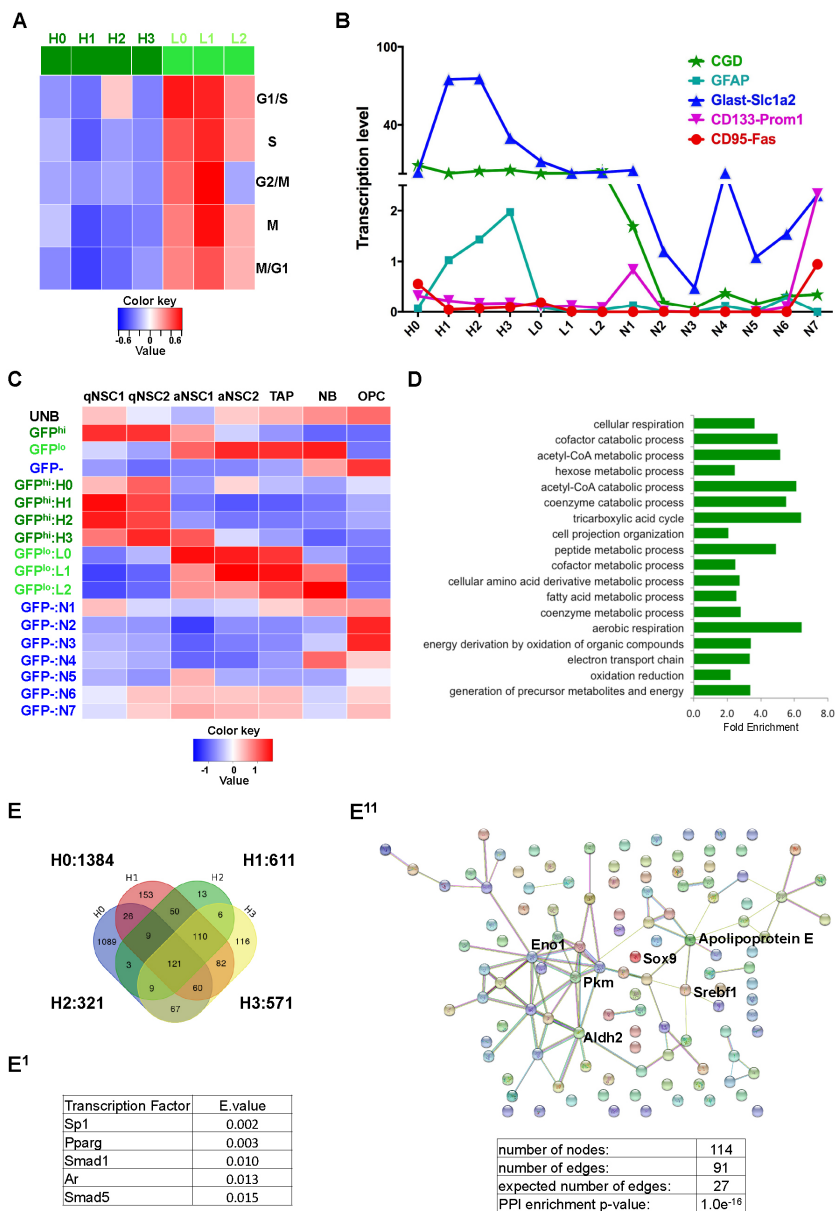


Fig. S4. This figure accompanies Figure 4. Granular view provided by analysis of the fourteen groups in the adult SVZ. (A) Cell cycle signature analysis of the GFP<sup>hi</sup> and GFP<sup>lo</sup> (H&L) groups reveals heterogeneity. Subgroup L1 shows the highest mitotic index for all seven signatures (see Fig. 2E). (B) Transcription levels of the five NSC markers within the fourteen populations (all standard errors < 1.2, units are normalized values scaled to 10,000 counts/cell). (C) Published single cell gene signatures fail to distinguish between the four GFP<sup>hi</sup> subgroups. (D) Known and unknown biological processes related to quiescent NSCs are revealed by a comprehensive list of genes derived from H0-H3 groups. (E-E<sup>11</sup>) An overlap analysis of genes commonly expressed in H0-H3 groups results in 121 genes (E), putative transcription factors associated to this list (TfactS analysis) (E<sup>1</sup>), and enrichment for a gene network is revealed (STRING analysis) (E<sup>11</sup>).

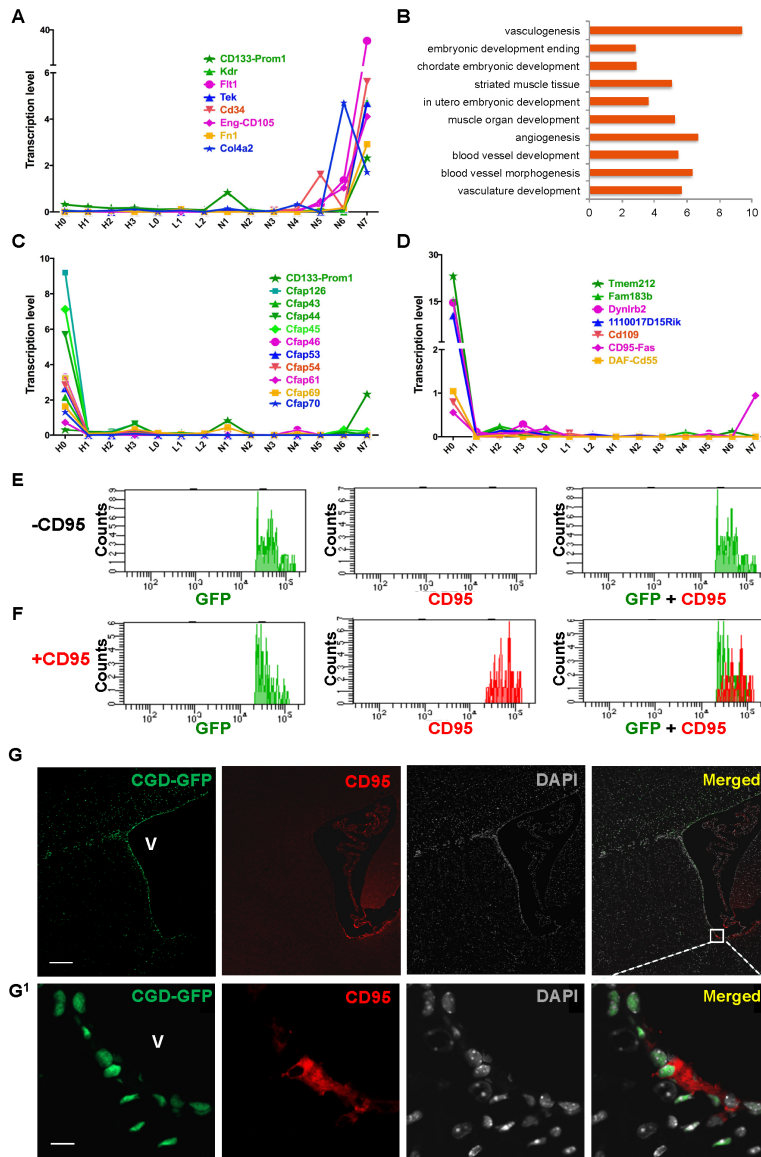


Fig. S5. This figure accompanies Figure 5. H0 subgroup gene analysis. (A) A gene network proposed as ependymal-like NSCs is enriched in the GFP<sup>-</sup>:N7 subgroup that by our analysis represents endothelial progenitor cells (all standard errors < 0.2, units are normalized values scaled to 10,000 counts/cell). (B) Gene ontology analysis of GFP<sup>-</sup>:N7 associate to vasculature and blood vessel development. (C) Cilia associated genes are uniquely expressed in GFP<sup>hi</sup>:H0 subgroup (all standard errors < 0.2, units are normalized values scaled to 10,000 counts/cell). (D) Expression profiles of seven potential markers for H0 subgroup cells (all standard errors < 2.0, units are normalized values scaled to 10,000 counts/cell). (E) and (F) Among the GFP<sup>+</sup> subgroups of cells, CD95 antibody preferentially labels cells with higher levels of GFP proteins. (G) IHC analysis reveals some CD95<sup>+</sup> cells co-localized with CGD-GFP at ventral SVZ. Scale bars: 200µm (G) and 10µm (G<sup>1</sup>). V: lateral ventricle in (G and G<sup>1</sup>).

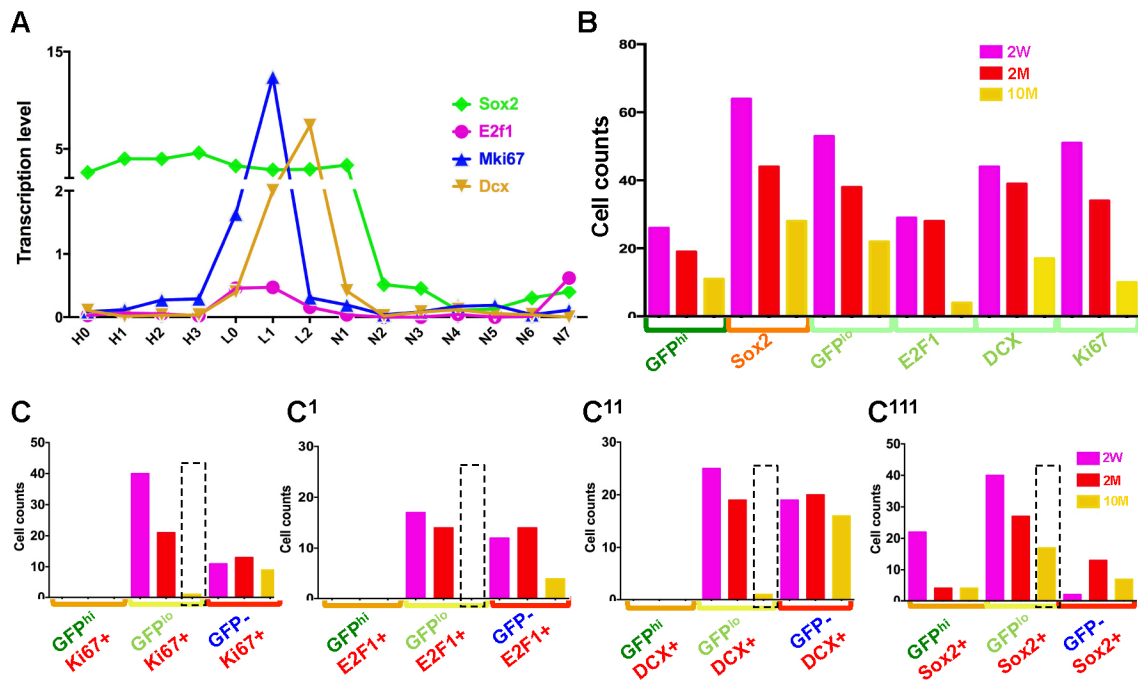


Fig. S6. This figure accompanies Figure 6. Age-related alteration of GFP<sup>hi</sup> and GFP<sup>lo</sup> lineage related genes. (A) Candidate genes for GFP<sup>hi</sup> qNSC and GFP<sup>lo</sup> progenitor cells illustrated in the adult SVZ analysis (all standard errors < 0.6, units are normalized values scaled to 10,000 counts/cell). (B) Concordant with age related GFP expression decrease, cells expressing four candidate genes decrease over age. (C-C<sup>111</sup>) The CGD-GFP-associated expression of the three candidate genes for NSC activation (C-C<sup>11</sup>, dashed line) drops more compared to the pan-NSC marker Sox2 (C<sup>111</sup>, dashed line).

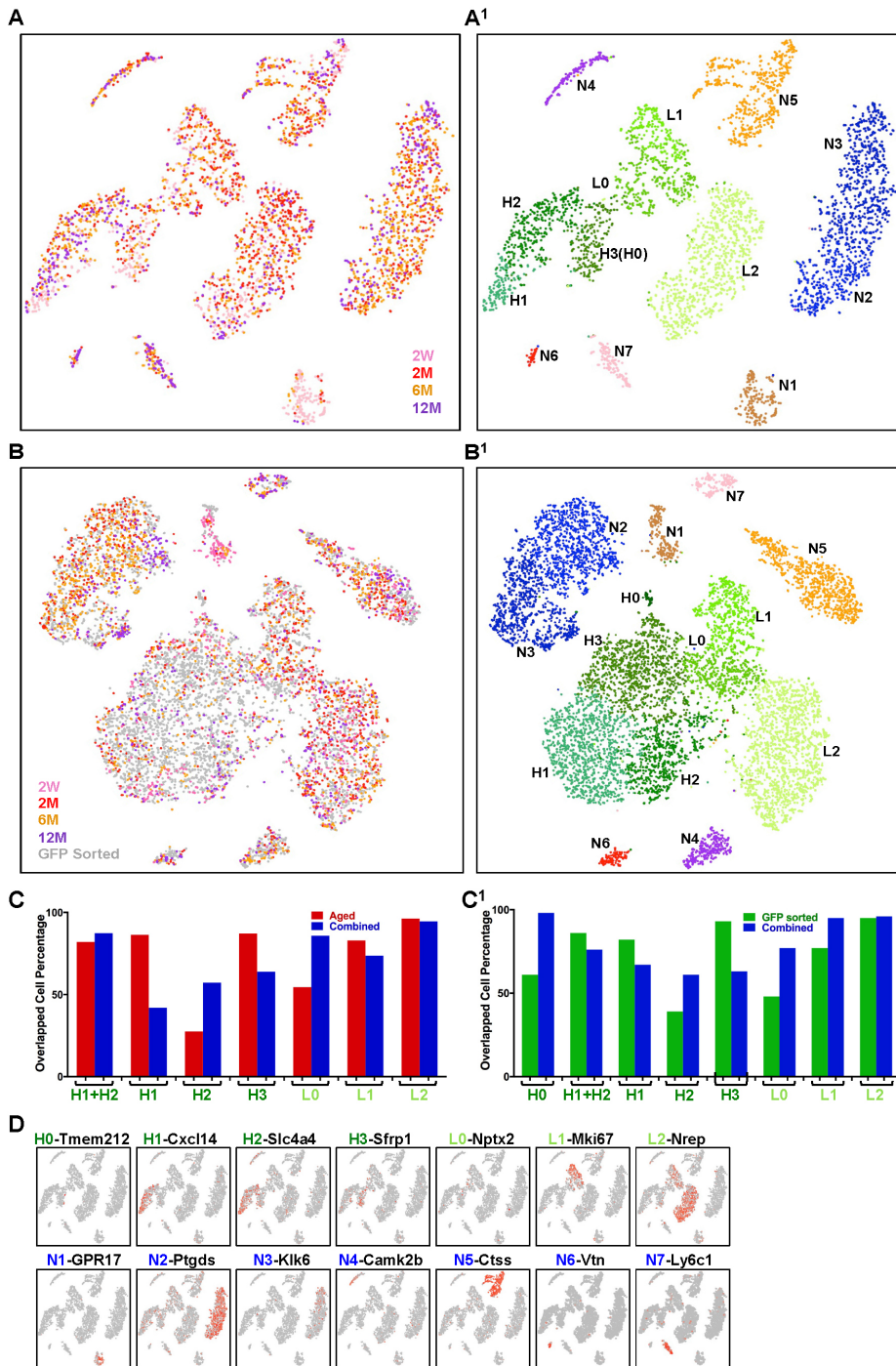


Fig. S7. This figure accompanies Figure 7. tSNE projection and marker gene analysis of the SVZ cells over age. (A-A') tSNE projection of the 5,600 single cells from two-week, two-month, six-month, and twelve-month SVZs reveals sample distribution (A) and thirteen groups (A'). The groups were assigned by comparing the cell constituents of the GFP sorted cohort (Fig. 3A), the Aged sample cohort (Fig. S7A'), and the combined analysis (Fig. S7B'). They were further confirmed by comparing unique genes from this analysis (Table S3C) with the unique genes from the GFP sorted cohort analysis (Table

S2C). DEGs were used instead when there were not enough unique genes (H2). (B-B<sup>1</sup>) tSNE projection of the comprehensive populations of the combined GFP sorted cohort with the Aged sample cohort (12,200 cells) demonstrates sample distribution (B) and fourteen conserved populations (B<sup>1</sup>). (C-C<sup>1</sup>) Each stem/progenitor group identity in the Aged sample cohort is examined further by comparing the cell constituents in (C) the combined versus the Aged samples and (C<sup>1</sup>) the combined versus the GFP sorted cohort. More than half of the cells are overlapped between the two cohorts in most cases. A significant amount of H2 cells are assigned into the H1 group in the combined analysis, possibly due in part to batch effect correction. (D) Marker genes identified in the two-month SVZ samples guide the identification of Seurat based clusters of SVZ cells from different ages. H0 cells (9 cells in total) are too few to form a group in the Aged cohort by themselves.

## **Movies S1 Legend**

To generate movie S1, the *CGD-GFP<sup>hi</sup>* cells were sorted from freshly dissociated mouse SVZ, incubated in Cytation 5 (BioTek), and imaged every hour for GFP and cell morphology.

## **Dataset Table Legends**

**Dataset Table S1, related to Figure 2E-2F and S2.** Transcriptomes (TS1A), differentiated expressed genes (TS1B), and GO analysis (TS1C&D) of the bulk RNAseq data derived from the *CGD-GFP<sup>hi</sup>* vs *GFP<sup>lo</sup>* samples. Complete lists of NSC genes derived from various studies including the *CGD-GFP<sup>hi</sup>* vs *GFP<sup>lo</sup>* samples are summarized in TS1E.

**Dataset Table S2, related to Figure 3, S3, 4, S4, and S5A-D.** Transcriptomes (TS2A), differentiated expressed genes (TS2B), unique genes (TS2C), and GO analysis (TS2D) of the fourteen cell groups derived from the two-month old murine SVZs. The 1914 genes derived from H0-H3 groups and associated GOs in Fig. S4D are summarized in (TS2E). The 121 overlapped G genes used in Fig. 4C&S4E, and related transcription factors are shown in (TS2F&G).

**Dataset Table S3, related to Figure 7 and S7.** Transcriptomes (TS3A), differentiated expressed genes (TS3B), and unique genes (TS3C) of the thirteen cell groups derived from the SVZs of two-week, two-month, six-month, and twelve-month mice. Transcriptomes of each cell group in the four aged SVZs are summarized in (TS3D). In lineage H1-H2, genes and GOs up-regulated in either two-month or twelve-month are shown in (TS3E&F). Transcription factors associated to the age-related change are indicated in (TS3G). In lineage H3-L1, genes and GOs up-regulated in either two-month



or twelve-month are shown in (TS3H&I). Potential transcription factors involved in the age-related blockage are summarized in (TS3J).

## SI References

1. E. Z. Macosko *et al.*, Highly Parallel Genome-wide Expression Profiling of Individual Cells Using Nanoliter Droplets. *Cell* **161**, 1202-1214 (2015).
2. W. Huang da, B. T. Sherman, R. A. Lempicki, Systematic and integrative analysis of large gene lists using DAVID bioinformatics resources. *Nat Protoc* **4**, 44-57 (2009).
3. P. Angerer *et al.*, destiny: diffusion maps for large-scale single-cell data in R. *Bioinformatics* **32**, 1241-1243 (2016).
4. E. S. Lein *et al.*, Genome-wide atlas of gene expression in the adult mouse brain. *Nature* **445**, 168-176 (2007).
5. H. Heberle, G. V. Meirelles, F. R. da Silva, G. P. Telles, R. Minghim, InteractiVenn: a web-based tool for the analysis of sets through Venn diagrams. *BMC Bioinformatics* **16**, 169 (2015).
6. A. Essaghir *et al.*, Transcription factor regulation can be accurately predicted from the presence of target gene signatures in microarray gene expression data. *Nucleic Acids Res* **38**, e120 (2010).
7. D. Szklarczyk *et al.*, The STRING database in 2017: quality-controlled protein-protein association networks, made broadly accessible. *Nucleic Acids Res* **45**, D362-D368 (2017).
8. X. Wang, A. Spandidos, H. Wang, B. Seed, PrimerBank: a PCR primer database for quantitative gene expression analysis, 2012 update. *Nucleic Acids Res* **40**, D1144-1149 (2012).
9. N. G. Sabbagha *et al.*, Alternative splicing in Acad8 resulting a mitochondrial defect and progressive hepatic steatosis in mice. *Pediatr Res* **70**, 31-36 (2011).
10. P. V. Kharchenko, L. Silberstein, D. T. Scadden, Bayesian approach to single-cell differential expression analysis. *Nat Methods* **11**, 740-742 (2014).