# Supplemental information
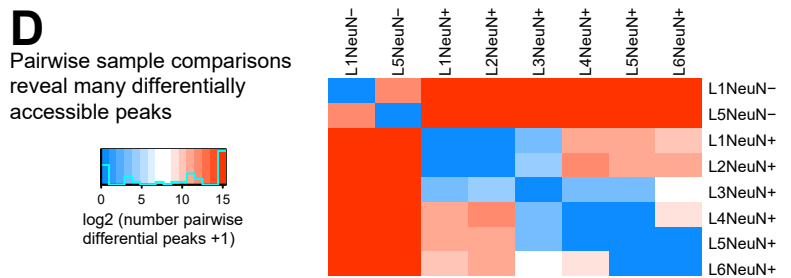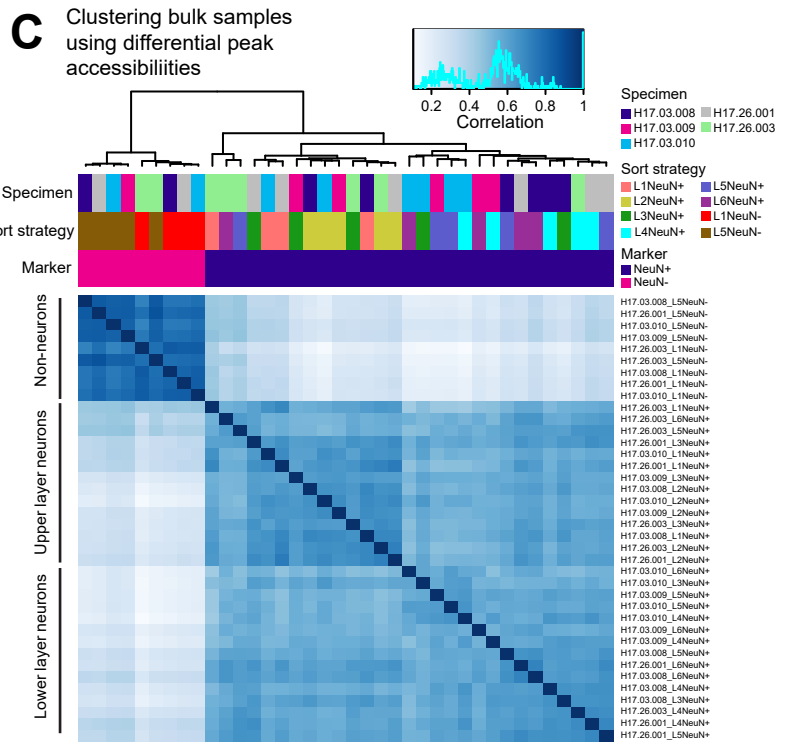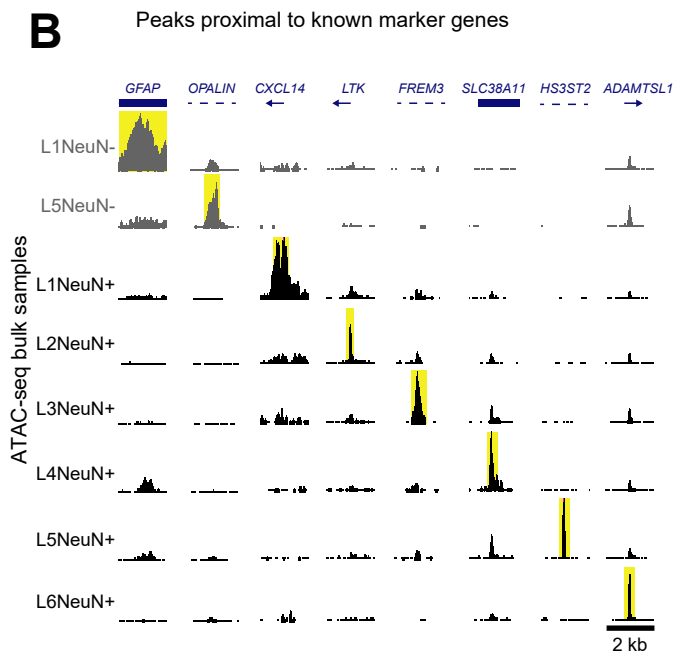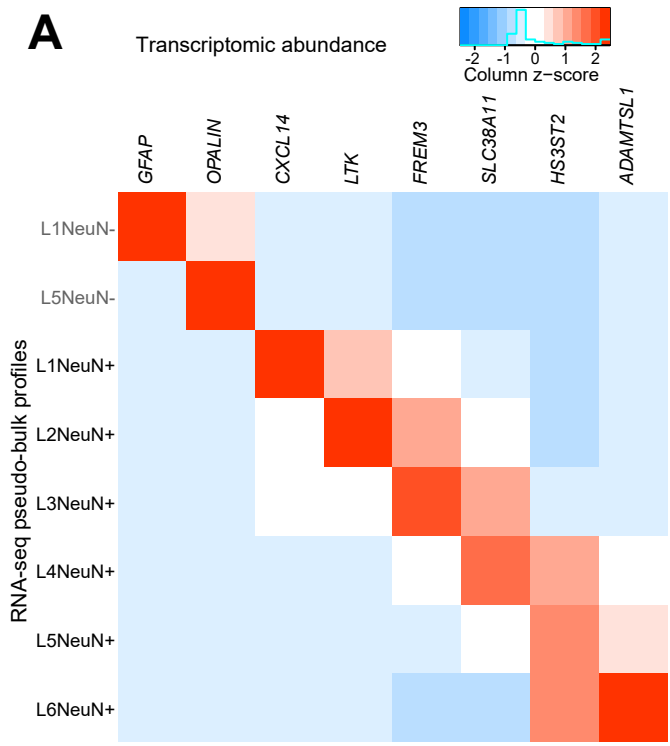
# Functional enhancer elements

# drive subclass-selective expression

# from mouse to primate neocortex

**John K. Mich, Lucas T. Graybuck, Erik E. Hess, Joseph T. Mahoney, Yoshiko Kojima, Yi Ding, Saroja Somasundaram, Jeremy A. Miller, Brian E. Kalmbach, Cristina Radaelli, Bryan B. Gore, Natalie Weed, Victoria Omstead, Yemeserach Bishaw, Nadiya V. Shapovalova, Refugio A. Martinez, Olivia Fong, Shenqin Yao, Marty Mortrud, Peter Chong, Luke Loftus, Darren Bertagnolli, Jeff Goldy, Tamara Casper, Nick Dee, Ximena Opitz-Araya, Ali Cetin, Kimberly A. Smith, Ryder P. Gwinn, Charles Cobbs, Andrew L. Ko, Jeffrey G. Ojemann, C. Dirk Keene, Daniel L. Silbergeld, Susan M. Sunkin, Viviana Gradinaru, Gregory D. Horwitz, Hongkui Zeng, Bosiljka Tasic, Ed S. Lein, Jonathan T. Ting, and Boaz P. Levi**

**A** Transcriptomic abundance

**B** Peaks proximal to known marker genes

**C** Clustering bulk samples using differential peak accessibiliities

**D** Pairwise sample comparisons reveal many differentially accessible peaks

**E** Peaks found in novel parts of genome from differential peak analysis

**Supplementary Figure 1 (Related to Figure 1): Bulk ATAC-seq data demonstrates differentially accessible chromatin elements around known marker genes, and in novel genomic regions.**
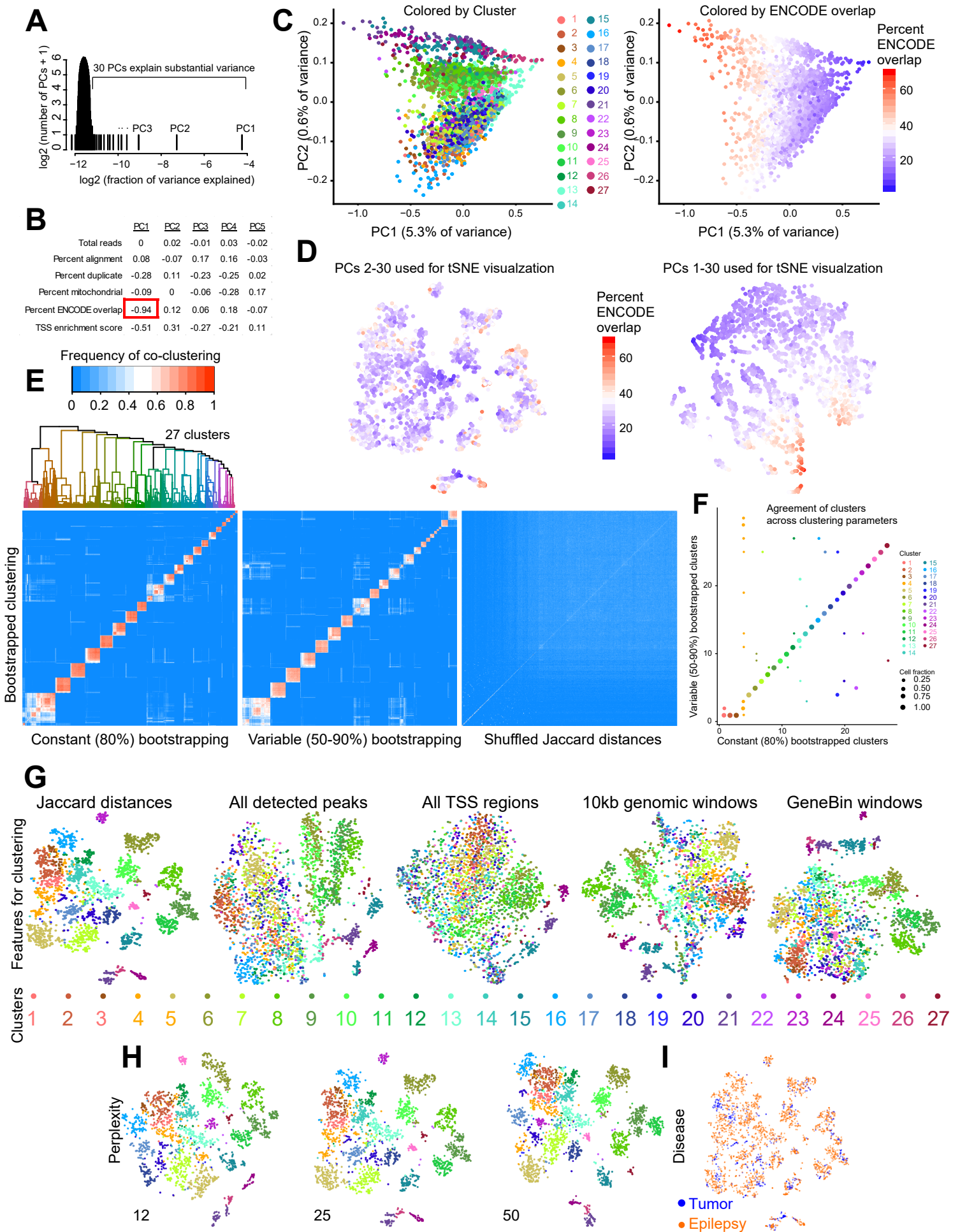
A) snRNA-seq data (Hodge et al., 2019), aggregated into pseudo-bulk profiles by weighted averages of gene CPM medians for 75 transcriptomic clusters. Weights were assigned by their frequencies within the eight sort strategies, and the heatmap is scaled by z-score within each column (gene). Relative expressions of eight sort strategy-specific marker genes are displayed.

B) Example sort strategy-specific peaks proximal to (<50kb distance to gene body) the eight sort strategy-specific transcriptomic marker genes. Pileups indicate aggregated data within a 2 kb genomic window across five independent experiments. In B and E, dashed lines indicate introns, thick lines indicate exons, and arrows indicate direction towards proximal marker gene. Yellow highlights demarcate sort strategy-specific chromatin accessibility peaks.

C) DiffBind (Ross-Innes et al., 2012) identification of 72,218 peaks that were differentially accessible among any pairwise comparison of sort strategies (FDR 0.01). Read counts within those 72,218 differentially accessible peaks then clustered samples using a correlation distance matrix, which revealed separate groupings of non-neuronal samples, and upper- and lower-layer neuronal samples. One sample was omitted from this analysis (H17.03.009 L1 NeuN+) because this sample appeared intermediate between NeuN+ and NeuN- cells, suggesting a failed sort.

D) Number of peaks differentiating each pairwise sample contrast.

E) Example sort strategy-specific peaks resulting from pairwise DiffBind differential peak analysis. These peaks were found in novel genomic regions (not proximal to known marker genes), and closest genes are shown.

**A**

30 PCs explain substantial variance

log2 (number of PCs + 1)

log2 (fraction of variance explained)

PC3  PC2  PC1

**B**

|  | PC1 | PC2 | PC3 | PC4 | PC5 |
|---|---|---|---|---|---|
| Total reads | 0 | 0.02 | -0.01 | 0.03 | -0.02 |
| Percent alignment | 0.08 | -0.07 | 0.17 | 0.16 | -0.03 |
| Percent duplicate | -0.28 | 0.11 | -0.23 | -0.25 | 0.02 |
| Percent mitochondrial | -0.09 | 0 | -0.06 | -0.28 | 0.17 |
| Percent ENCODE overlap | -0.94 | 0.12 | 0.06 | 0.18 | -0.07 |
| TSS enrichment score | -0.51 | 0.31 | -0.27 | -0.21 | 0.11 |

**C**

Colored by Cluster

PC2 (0.6% of variance)

PC1 (5.3% of variance)

1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27

Colored by ENCODE overlap

PC2 (0.6% of variance)

PC1 (5.3% of variance)

Percent ENCODE overlap

60 40 20

**D**

PCs 2-30 used for tSNE visualzation

PCs 1-30 used for tSNE visualzation

Percent ENCODE overlap

60 40 20

**E**

Frequency of co-clustering

0 0.2 0.4 0.6 0.8 1

27 clusters

Bootstrapped clustering

Constant (80%) bootstrapping

Variable (50-90%) bootstrapping

Shuffled Jaccard distances

**F**

Agreement of clusters across clustering parameters

Variable (50-90%) bootstrapped clusters

Constant (80%) bootstrapped clusters

Cluster
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27

Cell fraction
0.25 0.50 0.75 1.00

**G**

Features for clustering

Jaccard distances

All detected peaks

All TSS regions

10kb genomic windows

GeneBin windows

Clusters
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27

**H**

Perplexity

12

25

50

**I**

Disease

• Tumor
• Epilepsy

**Supplementary Figure 2 (Related to Figure 1): High confidence clustering for single nucleus ATAC-seq data.**

A) Histogram showing the percentages of variance explained by each principal component of the Jaccard single cell distance matrix. The first 30 principal components explain substantial variance within the dataset.

B) Correlation of the first five principal components with quality metrics. Principal component 1 was omitted from further analysis due to strong negative correlation with ENCODE overlap.

C) Single nuclei evaluated by principal component analysis, with nuclei colored by cluster membership (*left*). Three major groups of nuclei were separated by PC2. Single nuclei were also colored by ENCODE overlap percentage, which is strongly negatively correlated with PC1 (*right*).

D) tSNE plot to visualize either principal components 2 to 30 (*left*) or 1 to 30 (*right*). Note, PCs 2 to 30 permit clear groupings with no ENCODE overlap gradient, whereas PCs 1 to 30 result in blurred cluster separations with a gradient of ENCODE spanning the clusters.
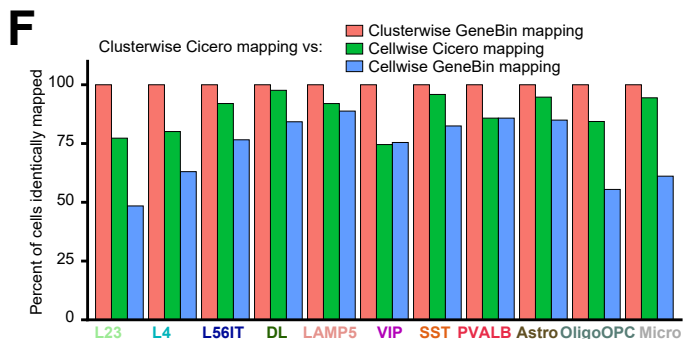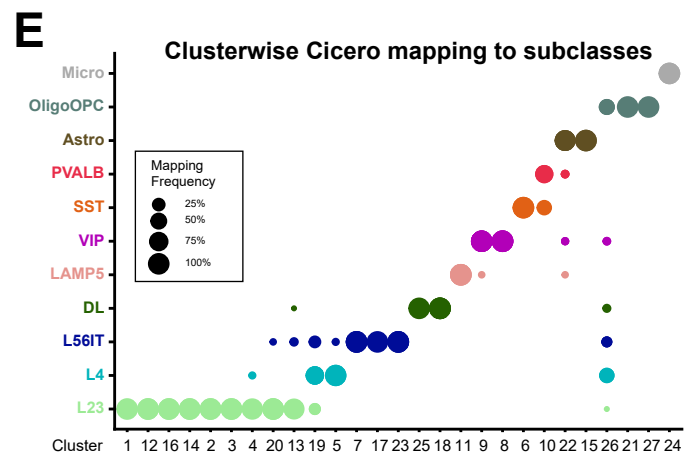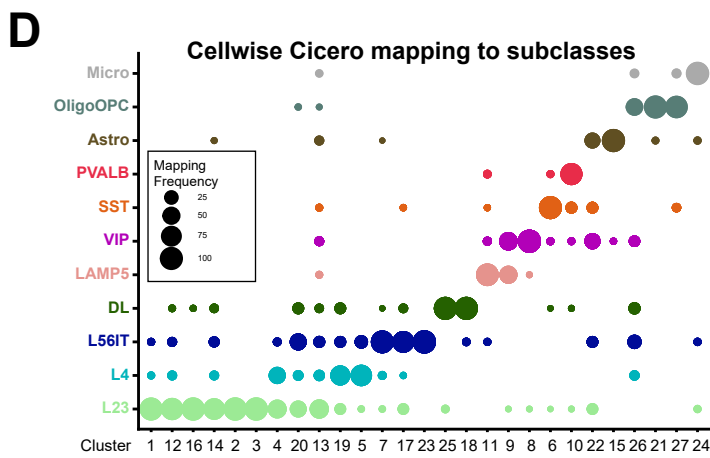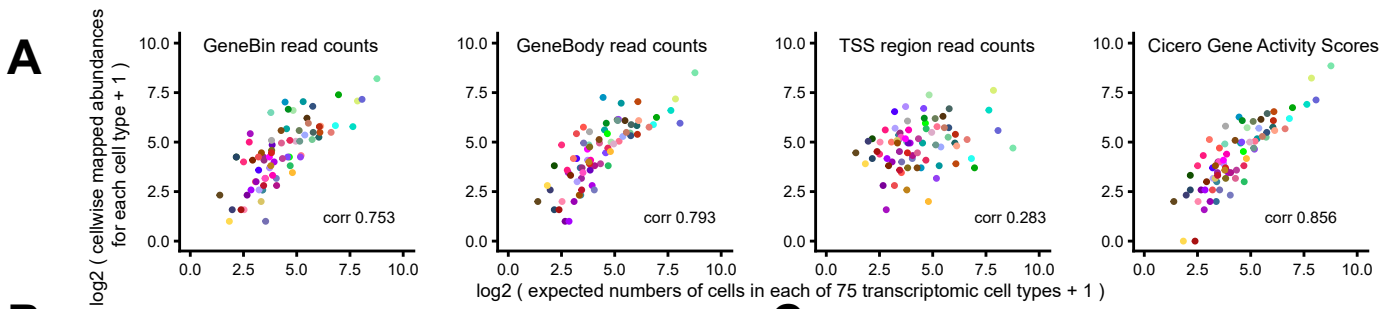
E) Bootstrapped iterative clustering to identify reproducible nuclear clusters. From the 2858 x 29 matrix of nuclei x principal component scores, we subsampled to either a constant 80% of nuclei (*left*) or a variable 50-90% of nuclei (*middle*), and calculated clusters using Jaccard-Louvain clustering (Tasic et al., 2018), which was repeated 200 times. Shuffled Jaccard distance matrix as input to PCA is shown (*right)*. Heatmaps display the frequency of co-clustering among nuclei. The constant 80% bootstrapping co-clustering matrix was used as input into Euclidean distance clustering, which yielded the final 27 clusters by cutting the tree to the major blocks of co-clustering nuclei. Nucleus order is not matched across the three plots.

F) Agreement between cluster memberships resulting from constant 80% bootstrapping and variable 50-90% bootstrapping, for most nuclei.

G) Visualization of nucleus groupings using five different feature sets (see Methods) using tSNE. Jaccard distances yielded clearest cell groupings. Cluster colors are applied in both (G) and (H).

H) Different perplexity parameters for tSNE visualization of cell groupings. Nuclear cluster groupings are evident at a wide range of perplexity values.

I) Visualization of disease status (tumor or epilepsy) for nuclei. Note that nuclei largely intermix regardless of disease status.

**Supplementary Figure 3 (Related to Figure 1): Mapping ATAC-seq clusters to RNA-seq cell types and subclasses.**

A) Expected abundances of each of the 75 transcriptomic cell types (Hodge et al., 2019), correlated with observed abundances of those cell types, using four different methods for computing gene-level information from each nucleus: *far left*: read counts in gene bins, *middle left*: read counts in gene bodies, *middle right:* read counts in 10kb-extended TSS regions, and *far right:* Cicero gene activity scores. Correlation values are Pearson correlation statistics between log-transformed expected and observed abundances plus one, for each of the 75 transcriptomic cell types. Computing gene-level information using cicero gene activity scores results in the greatest correlation between expectation and observation for cell type abundances.

B) Bootstrapped mapping of single nuclei ("cellwise") to 75 transcriptomic cell types. Dot sizes indicate the frequencies of cell type mappings within each of the 27 ATAC-seq clusters.
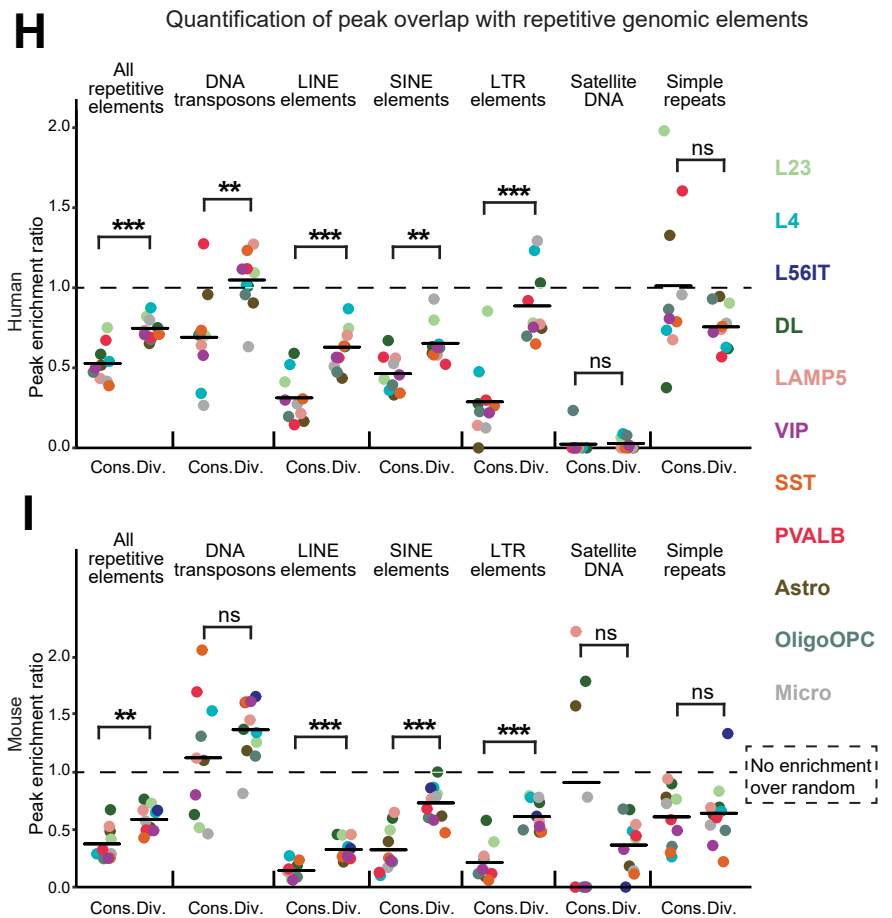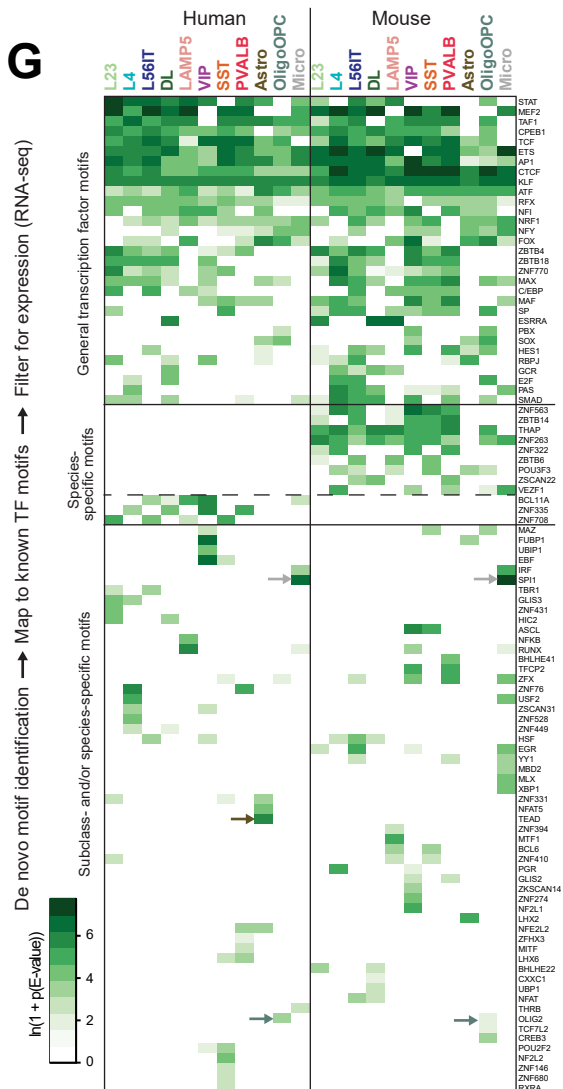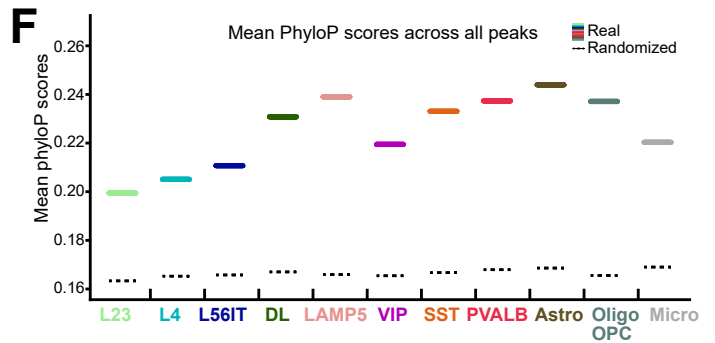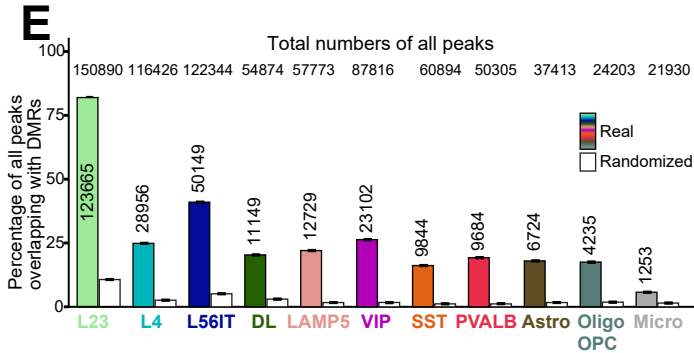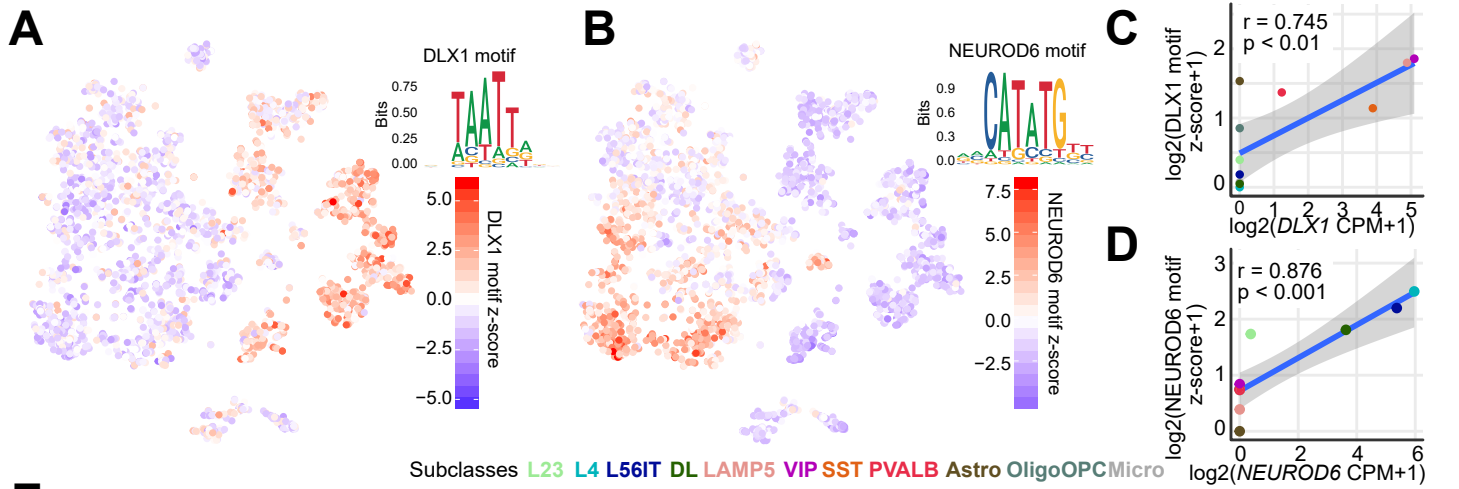
C) Bootstrapped mapping of clusters ("clusterwise") to 75 transcriptomic cell types. Dot sizes indicate the frequency each cluster maps to each transcriptomic cell type.

D) Bootstrapped mapping of single nuclei ("cellwise") to 11 transcriptomic cell type subclasses. Dot sizes indicate the frequencies of subclass mappings within each of the 27 ATAC-seq clusters.

E) Bootstrapped mapping of clusters ("clusterwise") to 11 transcriptomic cell type subclasses. Dot sizes indicate the frequency each cluster maps to each subclass. This plot represents the <u>final</u> mapped subclass assigned as the most frequent mapping for each cluster, which are used throughout the text.

F) Correlation of subclass mappings for all cells using four different mapping techniques. Overall, most cells are identically mapped to the same subclass with most of the techniques, with especially good agreement between both clusterwise mapping techniques.

G) Correlation between RNA-seq and ATAC-seq dataset layerwise distributions for the 11 subclasses. Most of the subclasses are observed in similar layer distributions in both datasets.

**A** DLX1 motif

**B** NEUROD6 motif

**C** r = 0.745, p < 0.01

**D** r = 0.876, p < 0.001

Subclasses L23 L4 L56IT DL LAMP5 VIP SST PVALB Astro OligoOPC Micro

**E** Total numbers of all peaks

**F** Mean PhyloP scores across all peaks

**G** Human Mouse

De novo motif identification → Map to known TF motifs → Filter for expression (RNA-seq)

General transcription factor motifs

Species-specific motifs

Subclass- and/or species-specific motifs

ln(1 + p[E-value])

**H** Quantification of peak overlap with repetitive genomic elements

Human — Peak enrichment ratio

All repetitive elements / DNA transposons / LINE elements / SINE elements / LTR elements / Satellite DNA / Simple repeats

Cons. Div.

**I** Mouse — Peak enrichment ratio

No enrichment over random

**Supplementary Figure 4 (Related to Figure 2): Properties of human neocortical cell subclass-specific accessible genomic elements.**

A-B) Nuclei visualized by tSNE and colored by motif accessibilities for A) DLX1 and B) NEUROD6 as calculated by chromVAR (Schep et al., 2017). *DLX1* transcripts are specifically detected in inhibitory neurons (Hodge et al., 2019).

C-D) Correlation between motif accessibilities and transcript abundances across cell subclasses for C) DLX1 and D) NEUROD6 (grouping by average for motif accessibility, and by sum for transcript abundances). r, Pearson correlation coefficient. Two-tailed paired t-tests for significant correlation: DLX1 $t$ = 3.0 df = 9 p < 0.01; NEUROD6 $t$ = 5.4 df = 9 p < 0.001.
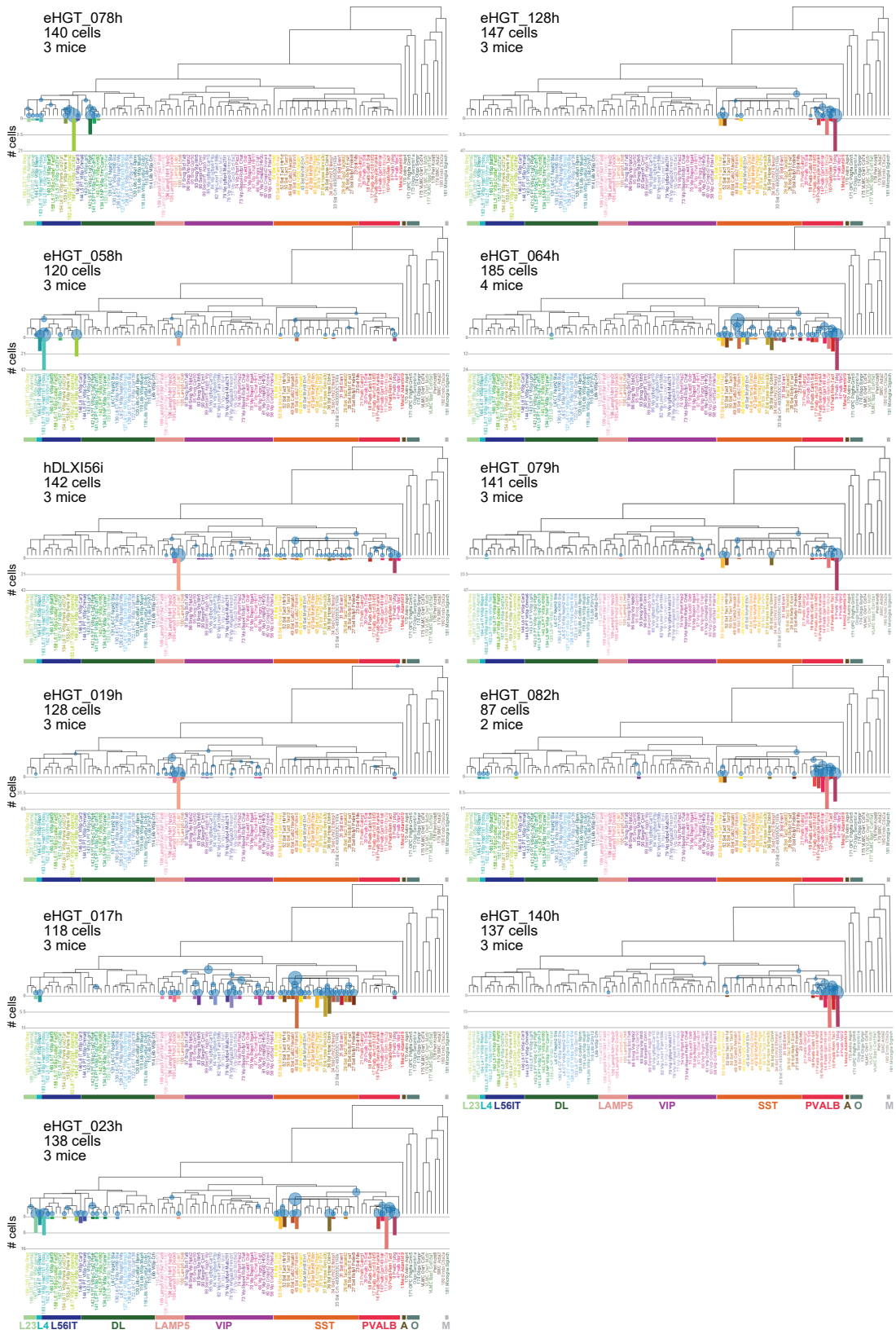
E) Percent overlap of ATAC-seq peaks with previously identified DMRs (Lister et al., 2013; Luo et al., 2017), comparing real peaks to randomized peak positions. Absolute numbers of detected peaks and peak-DMR overlaps are shown. Error bars represent standard deviation across 100 bootstrapped iterations using a subsampling rate of 80%.

F) Mean phyloP scores across all peaks for cell subclass ATAC-seq peaks (colored line), compared to randomized peak positions (broken line).

G) Active transcriptional regulators in human and mouse brain cell subclasses, revealed by motifs in ATAC-seq peaks and gene expression by transcriptomics (Tasic et al., 2018; Hodge et al., 2019). E-value indicates the p-value from Fisher's exact test, corrected for multiple testing as calculated by MEME-CHIP. Arrows indicate strong and specific microglial enrichments for SPI1/PU.1 (gray) and for TEAD in human astrocytes (brown) and for OLIG2 in oligodendrocytes/OPCs (cadet blue).
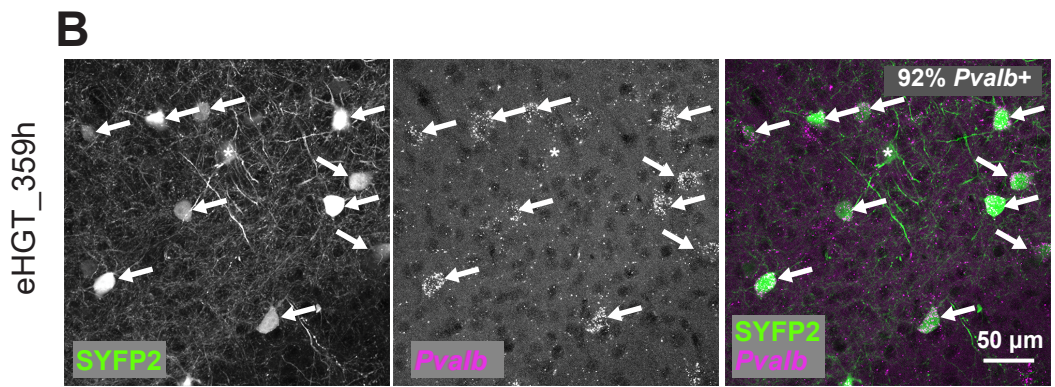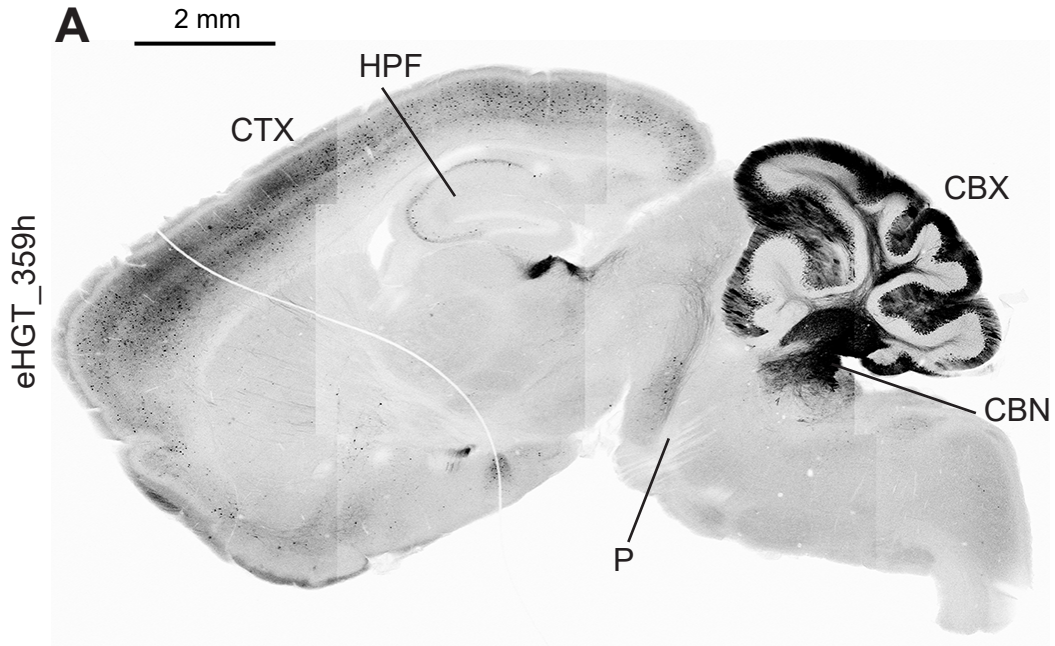
H-I) Overlap of conserved or divergent peaks by cell subclass with multiple classes of repetitive genomic elements in both human (H) and mouse (I) sn/scATACseq datasets. An enrichment value of 1.0 corresponds to no fold change between real and random peak enrichment. Black bars represent the mean across the eleven neocortical cell subclasses. Heteroscedastic t-tests: *** p < 0.001, ** p < 0.01, ns not significant. Human all elements t = 5.2, df = 14.5; human DNA

transposons t = 3.4, df = 14.8; human LINE t = 5.3, df = 18.3; human SINE t = 3.8, df = 18.8; human LTR t = 6.1, df = 18.3; human satellite t = 0.2, df = 12.4; human simple t = 1.6, df = 10.1; mouse all elements t = 3.7, df = 16.9; mouse DNA transposons t = 1.3, df = 12.7; mouse LINE t = 5.2, df = 18.5; mouse SINE t = 5.2, df = 18.0; mouse LTR t = 6.1, df = 17.5; mouse satellite t = 1.6, df = 9.7; mouse simple t = 0.3, df = 19.0.

eHGT_078h
140 cells
3 mice

eHGT_128h
147 cells
3 mice

eHGT_058h
120 cells
3 mice

eHGT_064h
185 cells
4 mice

hDLXl56i
142 cells
3 mice

eHGT_079h
141 cells
3 mice

eHGT_019h
128 cells
3 mice

eHGT_082h
87 cells
2 mice

eHGT_017h
118 cells
3 mice

eHGT_140h
137 cells
3 mice

eHGT_023h
138 cells
3 mice

L23 L4 L56IT    DL    LAMP5   VIP    SST    PVALB  A O    M

**Supplementary Figure 5 (Related to Figure 3): Cell type validation of enhancer-AAV-labeled cells via scRNA-seq.** Numbers of sorted labeled cells with each enhancer-AAV vector shown in Figures 3 and 4, mapped to the cell type transcriptomic taxonomy of mouse VISp (Tasic et al., 2018). Dendrogram leaves represent 111 transcriptomic cell types. Circles on the dendrogram represent the number of cells that could be mapped to that point in the dendrogram (starting from the root) and bar plots below the leaves represent the number of each cell type recovered that mapped to that final leaf.
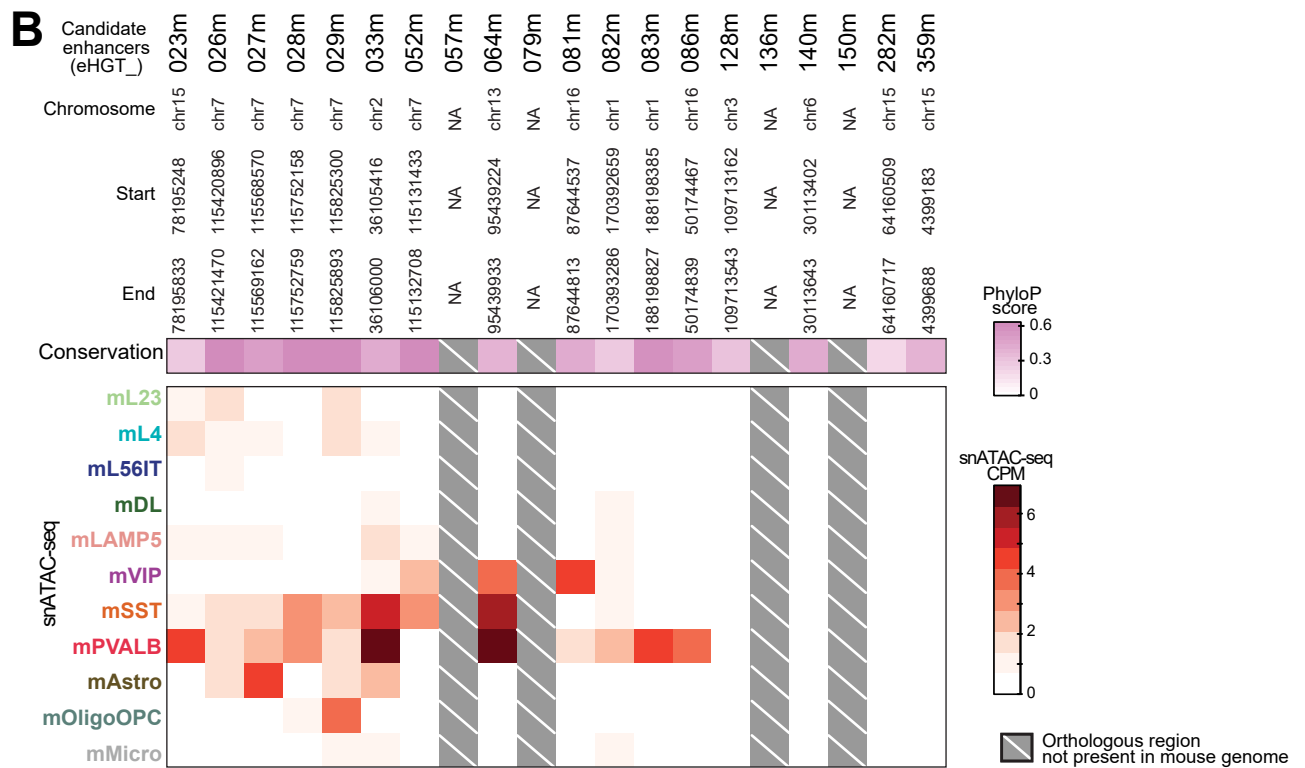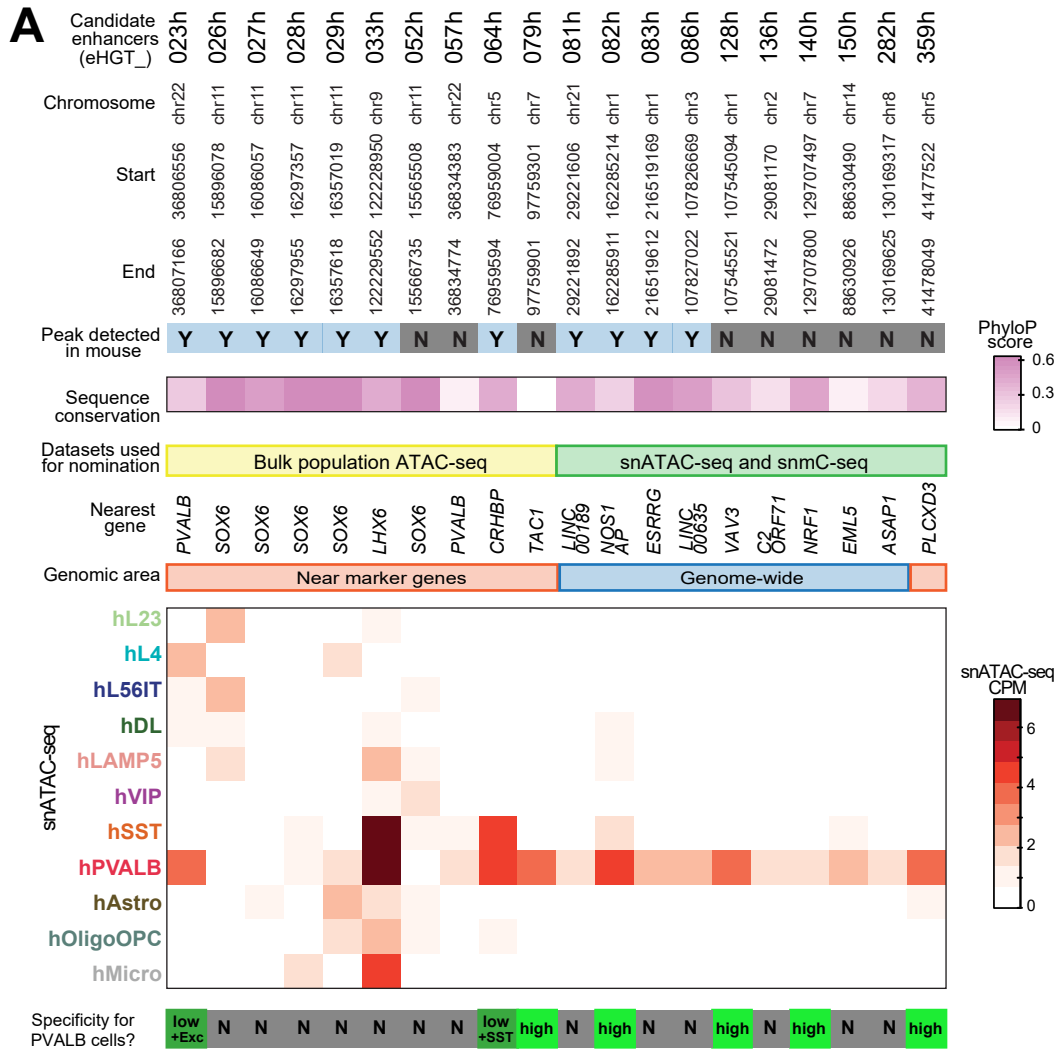
**A**

2 mm

CTX

HPF

CBX

CBN

P

eHGT_359h

**B**

eHGT_359h

SYFP2

*Pvalb*

92% *Pvalb+*

SYFP2
*Pvalb*

50 µm

**Supplementary Figure 6 (Related to Figure 4): eHGT_359h labels neocortical PVALB cells with high specificity, and also subcortical zones of *Pvalb* expression.**

A) Whole-brain labeling pattern by eHGT_359h. Abbreviations: CTX cerebral cortex, HPF hippocampal formation, P pons, CBX cerebellar cortex, CBN cerebellar nuclei.

B) mFISH in L2/3 of VISp demonstrating positive labeling of $Pvalb^+$ cells (arrows) by eHGT_359h. An asterisk denotes an off-target $Pvalb^-$ labeled cell. Data represents n = 1 experiment.

**Supplementary Figure 7 (Related to Figure 4): Accessibility of candidate PVALB enhancers in human and mouse neocortical subclasses.**

A) Twenty candidate PVALB enhancers from human epigenetic data characterized by conservation, method of identification, genomic location, and ATAC-seq profiles across cell subclasses. Coordinates correspond to hg38 genome in (A).

B) Orthologous mouse regions characterized by the same metrics. No mouse regions orthologous to eHGT_057h, 079h, 136h, and 150h could be identified, using liftOver with minMatch parameter set to either 0.6 or 0.5. CPM counts per million. Coordinates correspond to mm10 genome in (B).