

## Supplemental information

### Single nucleus multi-omics identifies

### human cortical cell regulatory genome diversity

Chongyuan Luo, Hanqing Liu, Fangming Xie, Ethan J. Armand, Kimberly Siletti, Trygve E. Bakken, Rongxin Fang, Wayne I. Doyle, Tim Stuart, Rebecca D. Hodge, Lijuan Hu, Bang-An Wang, Zhuzhu Zhang, Sebastian Preissl, Dong-Sung Lee, Jingtian Zhou, Sheng-Yong Niu, Rosa Castanon, Anna Bartlett, Angeline Rivkin, Xinxin Wang, Jacinta Lucero, Joseph R. Nery, David A. Davis, Deborah C. Mash, Rahul Satija, Jesse R. Dixon, Sten Linnarsson, Ed Lein, M. Margarita Behrens, Bing Ren, Eran A. Mukamel, and Joseph R. Ecker

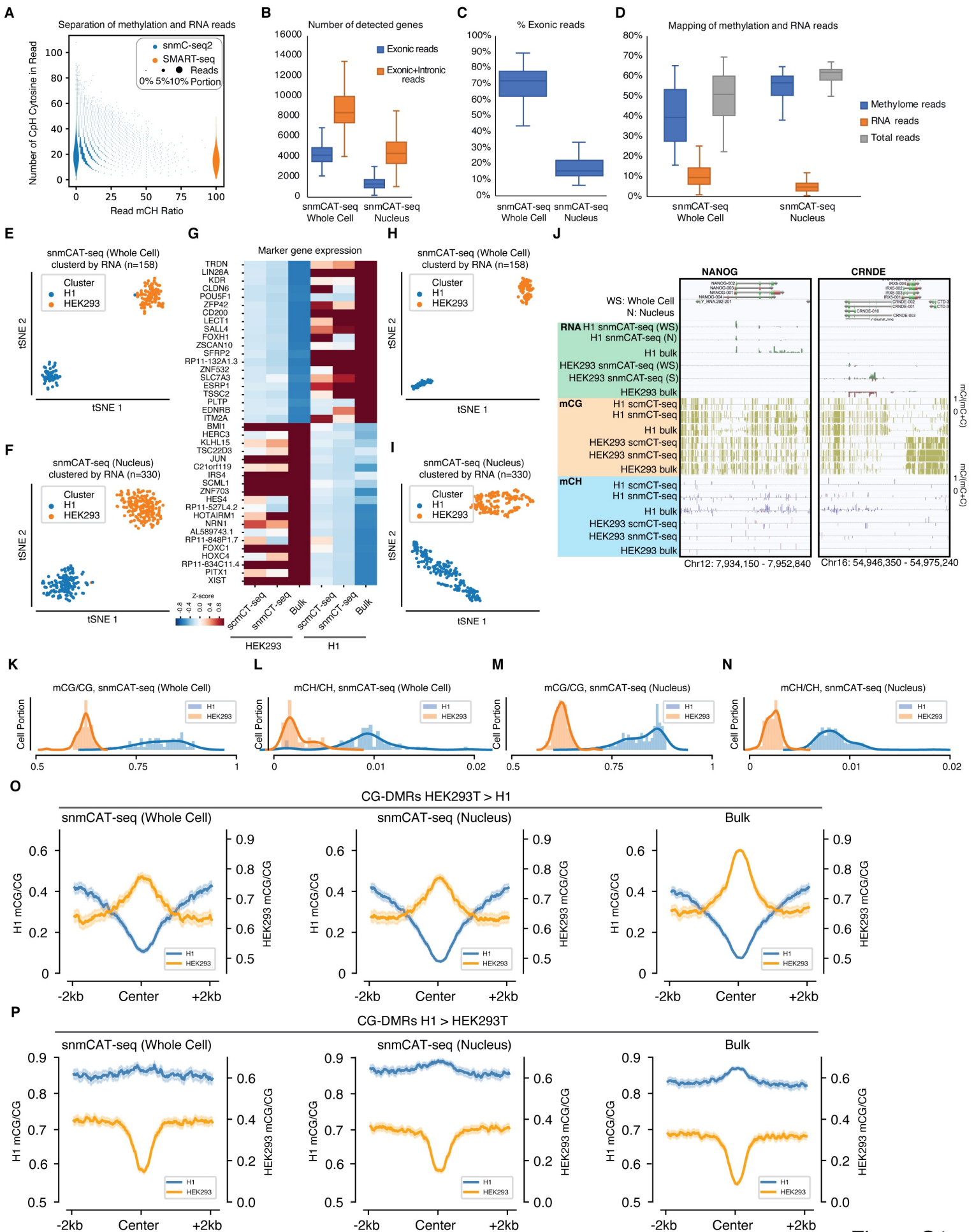


Figure S1

**Figure. S1 snmCAT-seq captures transcriptome and DNA methylation signatures of H1 & HEK293 cells, related to Figure 1.** (A) The specificity for classifying methylation (snmC-seq2) and transcriptome (snRNA-seq) reads plotted as a function of the number of CpH cytosine in the reads. (B-D) The number of detected genes (B), percentage of mapped reads that are located in exons (C), and mapping rates of methylation and RNA reads (D) for snmCAT-seq (Whole Cell) and snmCAT-seq (Nucleus). (E-F) Separation of H1 and HEK293T cells by tSNE using transcriptome reads extracted from snmCAT-seq (Whole Cell) (E) or snmCAT-seq (Nucleus) (F) datasets. (G) snmCAT-seq (Whole Cell) and snmCAT-seq (Nucleus) detect genes specifically expressed in H1 or HEK293T cells. (H-I) Separation of H1 and HEK293T cells by tSNE using DNA methylation information extracted from snmCAT-seq (Whole Cell) (H) or snmCAT-seq (Nucleus) (I) datasets. (J) Browser view of NANOG and CRNDE loci. (K-N) Distribution of mCG and mCH levels for single H1 and HEK293 cells/nuclei as determined by snmCAT-seq (Whole Cell) and snmCAT-seq (Nucleus). (O-P) snmCAT-seq (Whole Cell) and snmCAT-seq (Nucleus) recapitulate bulk mCG patterns at CG-DMRs showing greater mCG levels in HEK293T (O) or H1 (P) cells.

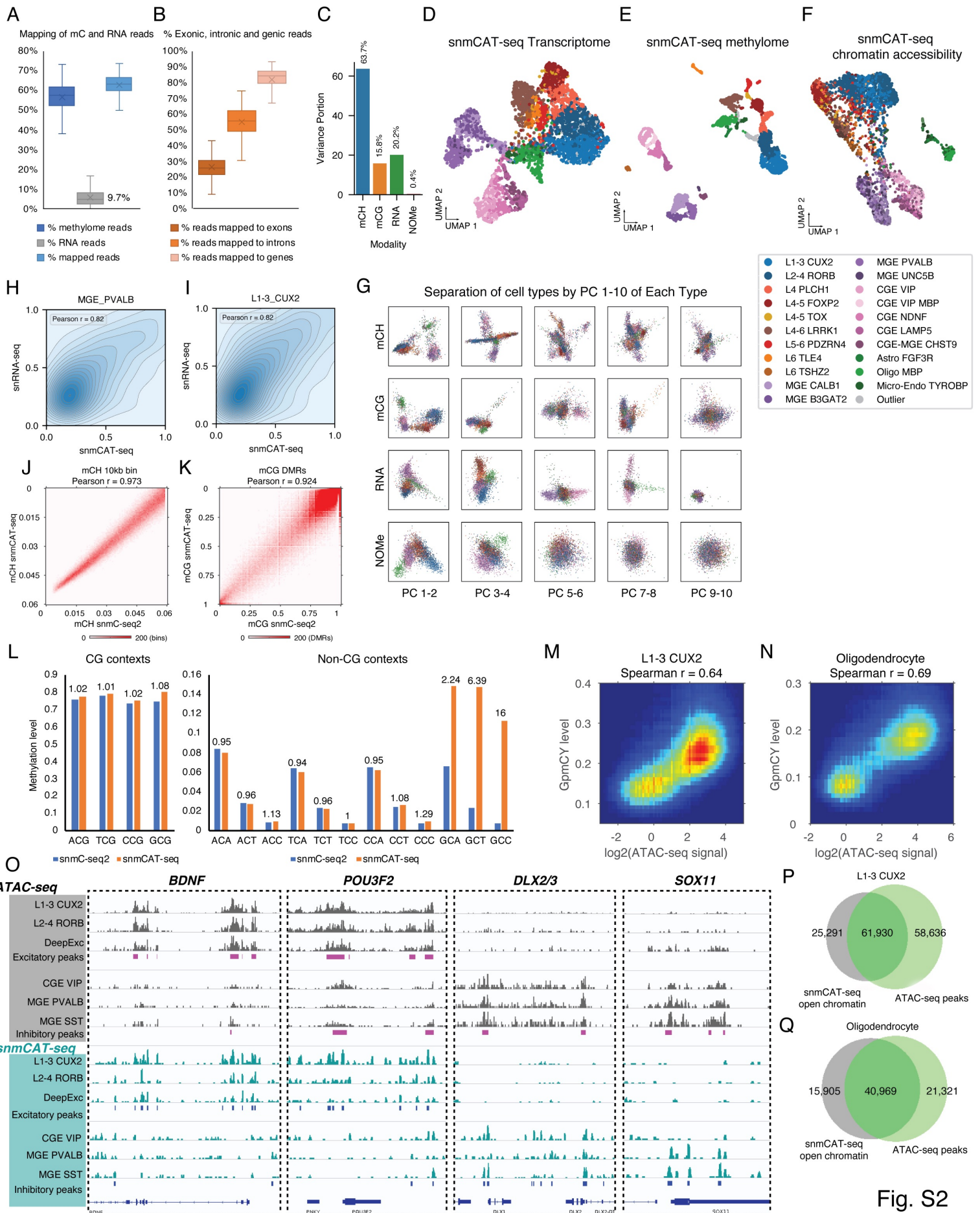


Fig. S2

**Figure S2. snmCAT-seq generates single-nucleus multi-omic profiles from human brain tissues, related to Figure 1.** (A) The fraction of total snmCAT-seq reads derived from methylome and transcriptome. (B) The fraction of snmCAT-seq transcriptome reads mapped to exons, introns or gene bodies. (C) The portion of variance explained by each data modality. (D-F) UMAP embedding of 4253 snmCAT-seq cells using single modality information: transcriptome (D), methylome (mCH and mCG, E) and chromatin accessibility (F). (G) The separation of cell types by the top 10 principal components (PCs) of each data type. (H-I) Pearson correlation of gene expression quantified by snmCAT-seq transcriptome and snRNA-seq in MGE PVALB (H) and L1-3 CUX2 (I) cells. (J) Pearson correlation of gene body non-CG methylation quantified with snmCAT-seq methylome and snmC-seq for MGE PVALB cells. (K) Pearson correlation of CG methylation at DMRs quantified with snmCAT-seq methylome and snmC-seq for MGE PVALB cells. (L) Genome-wide methylation level for all tri-nucleotide context (-1 to +2 position) surrounding cytosines shows the sequence specificity of GpC methyltransferase M.CviPI. (M-N) Spearman correlation between GCY methylation level and ATAC-seq signal at open chromatin sites in L1-3 CUX2 (M) and Oligodendrocyte (N) cells. (Q) Browser views of chromatin accessibility profiles generated by snmCAT-seq and snATAC-seq at cell-type-specific genes. (P-Q) Overlap of open chromatin peaks identified by snmCAT-seq and snATAC-seq in L1-3 CUX2 (P) and oligodendrocyte (Q) cells.

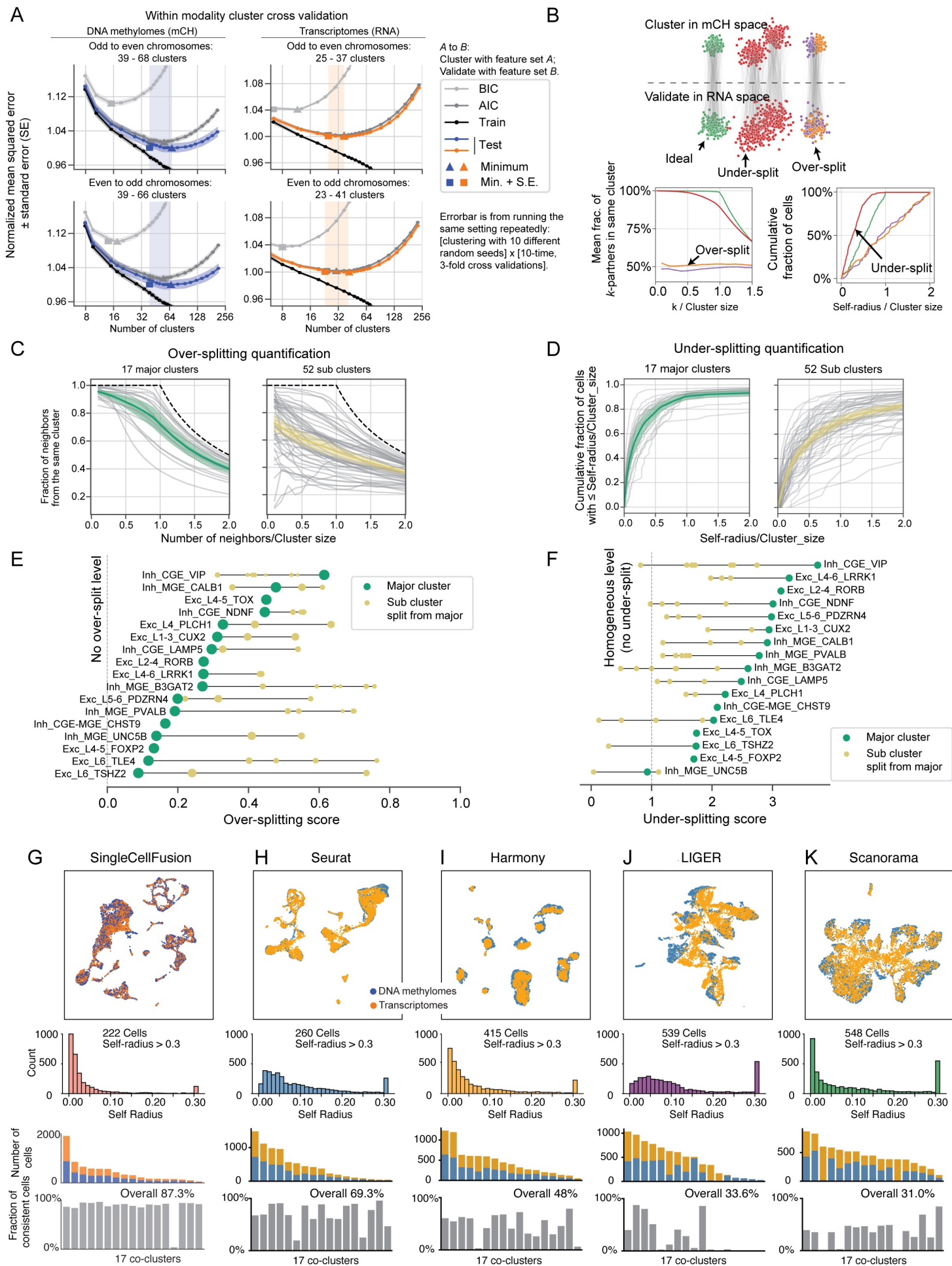


Figure S3

**Figure S3. Evaluation of cluster quality with paired transcriptome and methylome profiles, related to Figure 2.** (A) Intra-modality cross-validation of mCH- or RNA-defined clusters. Line plots show mean squared error between the single-cell profiles and cluster centroid, Akaike Information Criterion (AIC) and the Bayesian Information Criterion (BIC) as a function of the number of clusters. The shaded region in each sub-plot highlights the range between the minimum and the minimum + standard error for the curve of test-set error. For the analysis using snmCAT-seq mC information (left panels), gene body mCH profiles of odd (even) chromosomes were used for clustering whereas even (odd) chromosomes were used for testing. A similar analysis was performed using snmCAT-seq transcriptome information (right panels). (B) Schematic diagram of the over- and under-splitting analysis using matched single-cell methylome and transcriptome profiles, complementing Figure 2D. (C) The over-splitting quantification of mC-defined major clusters (n=17) and subclusters (n=52) was quantified by the fraction of cross-modal neighbors found in the same cluster defined by RNA. (D) The under-splitting of clusters was quantified as the cumulative distribution function of normalized self-radius for mC-defined major clusters and subclusters. For (C-D), gray lines represent individual clusters while colored lines represent means and confidence intervals. (E) Over-splitting score for each major cluster (in green) and associated sub-clusters (in yellow). Dot size of sub-clusters represents cluster size normalized by the size of their “mother” major cluster. (F) Under-splitting score for each major cluster (in green) and associated sub-clusters (in yellow). (G-K) Fusion of snmCAT-seq transcriptome and mC profiles using the Single Cell Fusion (G), Seurat (H), Harmony (I), LIGER (J), and Scanorama (K). For each computational data fusion

method, from top to bottom the first panel shows mC and RNA modalities on the joint UMAP embedding after data fusion. The second panel shows the normalized self-radius. The third and fourth panels show the co-cluster level cell composition and clustering accuracy.



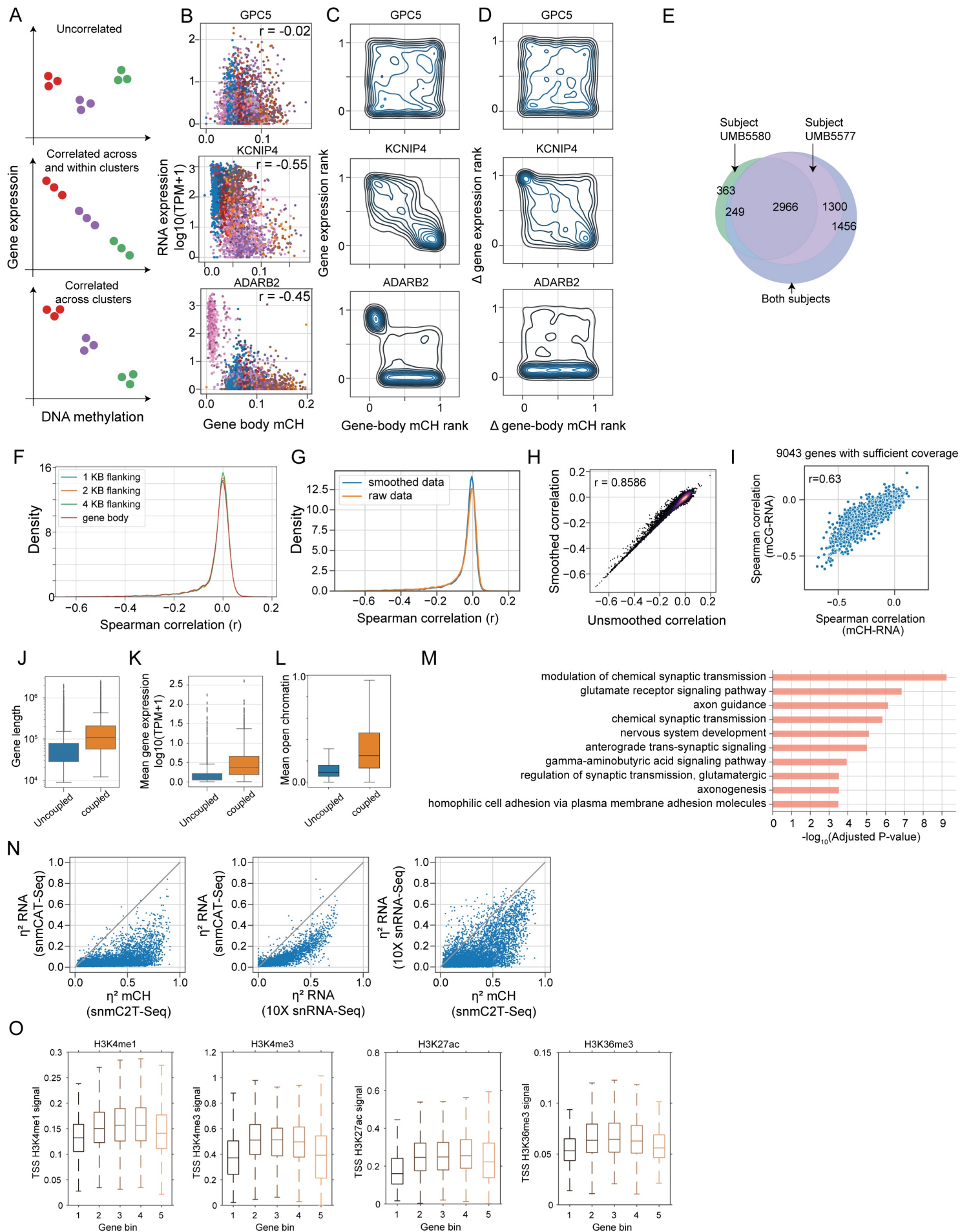
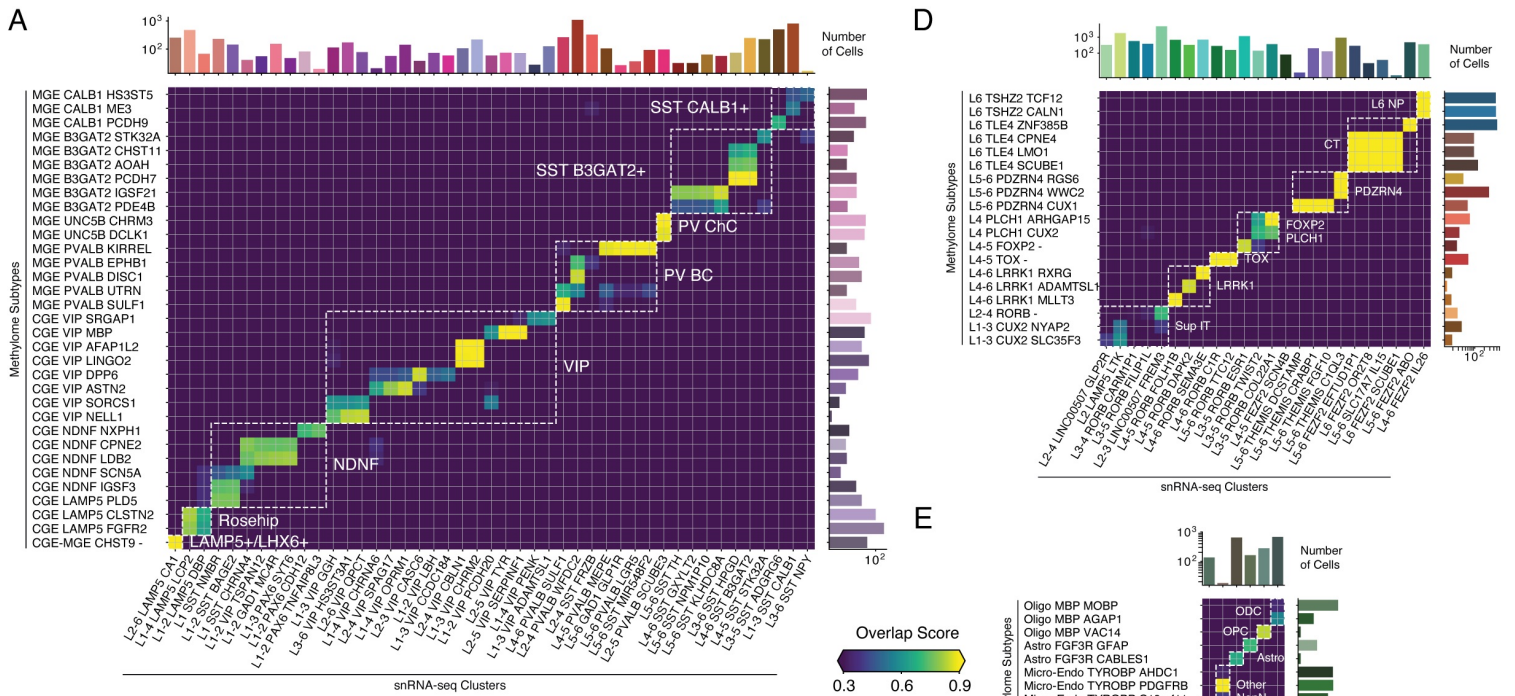


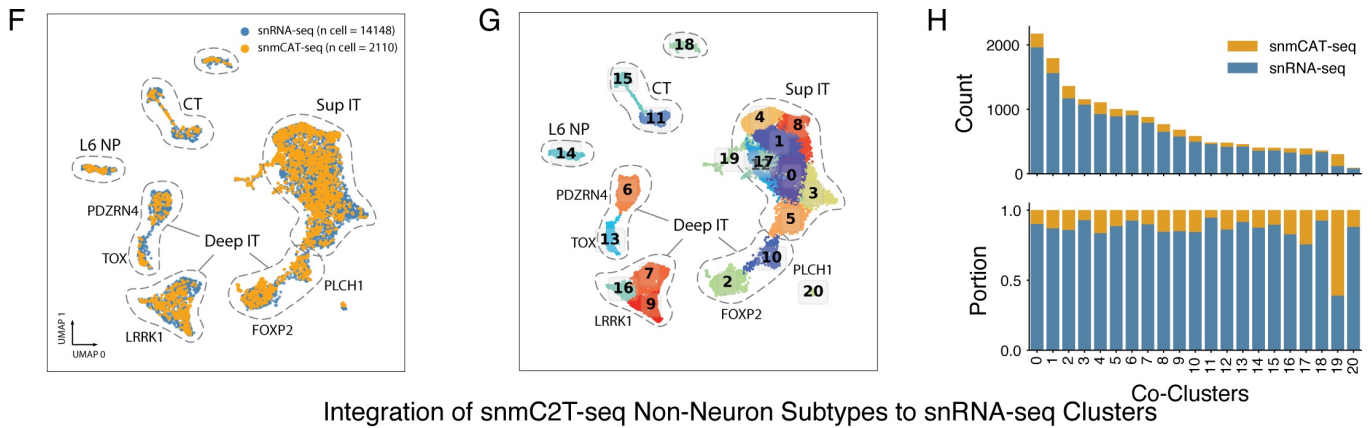
Figure S4

**Figure S4. Diverse correlations between gene expression and gene body mCH, related to Figure 3.** (A) Schematic diagram of 3 types of genes with different correlations between gene expression and DNA methylation. (B) Scatter plots of gene body mCH (unnormalized) and gene expression of example genes (KCNIP4, ADARB2, GPC5) across all neuronal cells. Cells are colored by major cell types defined in Figure 4. (C) Contour density plot of gene body mCH rank versus gene expression rank for 3 example genes: GPC5, KCNIP4, and ADARB2. (D) Contour density plot of delta gene body mCH rank (rank of gene body mCH - cluster mean gene body mCH) versus delta gene expression rank (rank of gene expression - cluster mean gene expression) for 3 example genes: GPC5, KCNIP4, and ADARB2. (E) Venn diagram showing a strong overlap of mCH-RNA coupled genes identified using cells from subject UMB5580, subject UMB5577, and from both subjects. (F) Distribution of Spearman correlation coefficient between gene expression and gene-level mCH quantified at gene body (red), gene body + 1 kilo-base upstream (blue), + 2 kilo-base upstream (orange) and + 4 kilo-base upstream (green). (G) Distribution of Spearman correlation coefficient between gene expression and gene-level mCH quantified at gene body with or without data smoothing before correlation analysis. (H) Scatter plot comparing the effect of data smoothing on Spearman correlation coefficients computed between gene expression and gene-level mCH. (I) Scatter plot comparing gene body mCH-RNA correlation versus gene body mCG-RNA correlation. (J-L) Boxplot of gene length (J), mean gene expression level (K) and mean open chromatin level (L) for mCH-RNA uncoupled and coupled genes. (M) Gene ontology enrichment of mCH-RNA coupled genes. (N) Scatter plots comparing the fraction of variance explained by cell type ( $\eta^2$ ) for each gene from

different datasets or data modalities: RNA (from snmCAT-seq), mCH (from snmCAT-seq) and 10X (snRNA-Seq from 10X protocols). (O) The distribution of different histone marks at TSS over 5 gene bins grouped according to gene expression ratio of early fetal (PCW 8-9) to adult (>2 yrs).



### Integration of snmCAT-seq Excitatory Subtypes to snRNA-seq Clusters



### Integration of snmC2T-seq Non-Neuron Subtypes to snRNA-seq Clusters

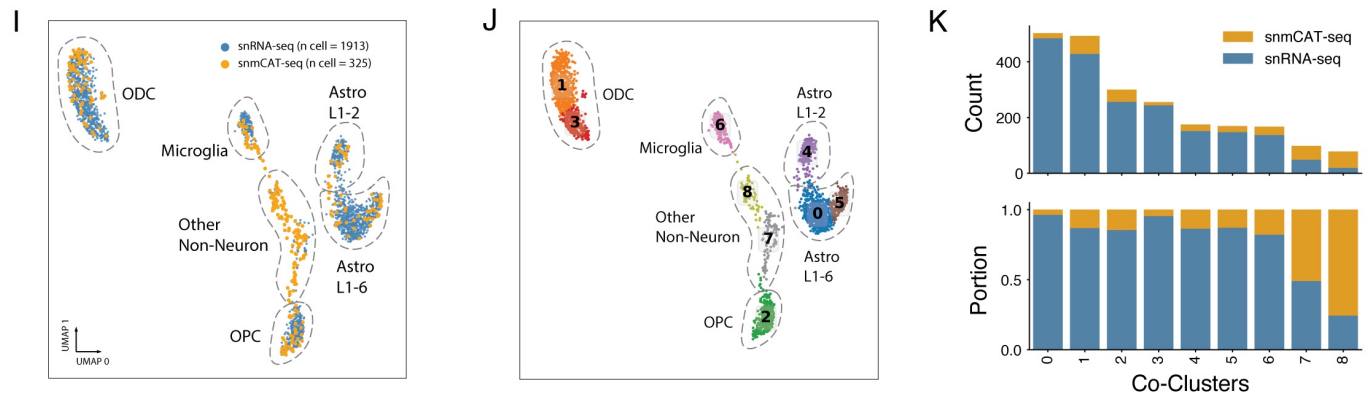


Figure S5

**Figure S5. snmCAT-seq recapitulates transcriptome and methylome signatures of neuronal subtypes, related to Figure 5.** (A) Confusion matrix showing the overlap scores between inhibitory subtypes identified by ensemble methylome analysis and snRNA-seq. Known inhibitory cell type groups are annotated by boxes. The upper bar plot indicates the snRNA-seq cell counts per cluster; the right bar plot indicates the snmCAT-seq cell counts per cluster. (B) Cluster 13 in Figure 5A-C includes snmCAT-seq transcriptome profiles with lower RNA read counts. (C) snmCAT-seq methylome profiles of single nuclei in cluster 13 (colored in blue) can be readily integrated with other inhibitory cells. (D-E) Confusion matrix showing the overlap scores between methylome ensemble subtypes and the snRNA-seq clusters for excitatory neuron clusters (D) and non-neuron clusters (E). Known cell type groups are annotated by boxes. The upper bar plot annotates the snRNA-seq cell counts per cluster, the right bar plot annotates the snmCAT-seq cell counts per cluster. (F-G) UMAP embedding of all excitatory neurons profiled by snmCAT-seq and snRNA-seq after MNN-based integration, colored by technology (F) and joint clusters (G). Known cluster groups are also circled and annotated on UMAP. (H) The composition of cells profiled by snmCAT-seq and snATAC-seq in excitatory neuron joint clusters. The upper and lower bar plots show the counts and portion of cells profiled by the two technologies in each joint cluster, respectively. (I-J) UMAP embedding of all non-neuronal cells from snmCAT-seq and snRNA-seq after integration, colored by technology (I) and joint clusters (J). (K) The composition of cells profiled by snmCAT-seq and snATAC-seq in non-neuronal cell joint clusters.

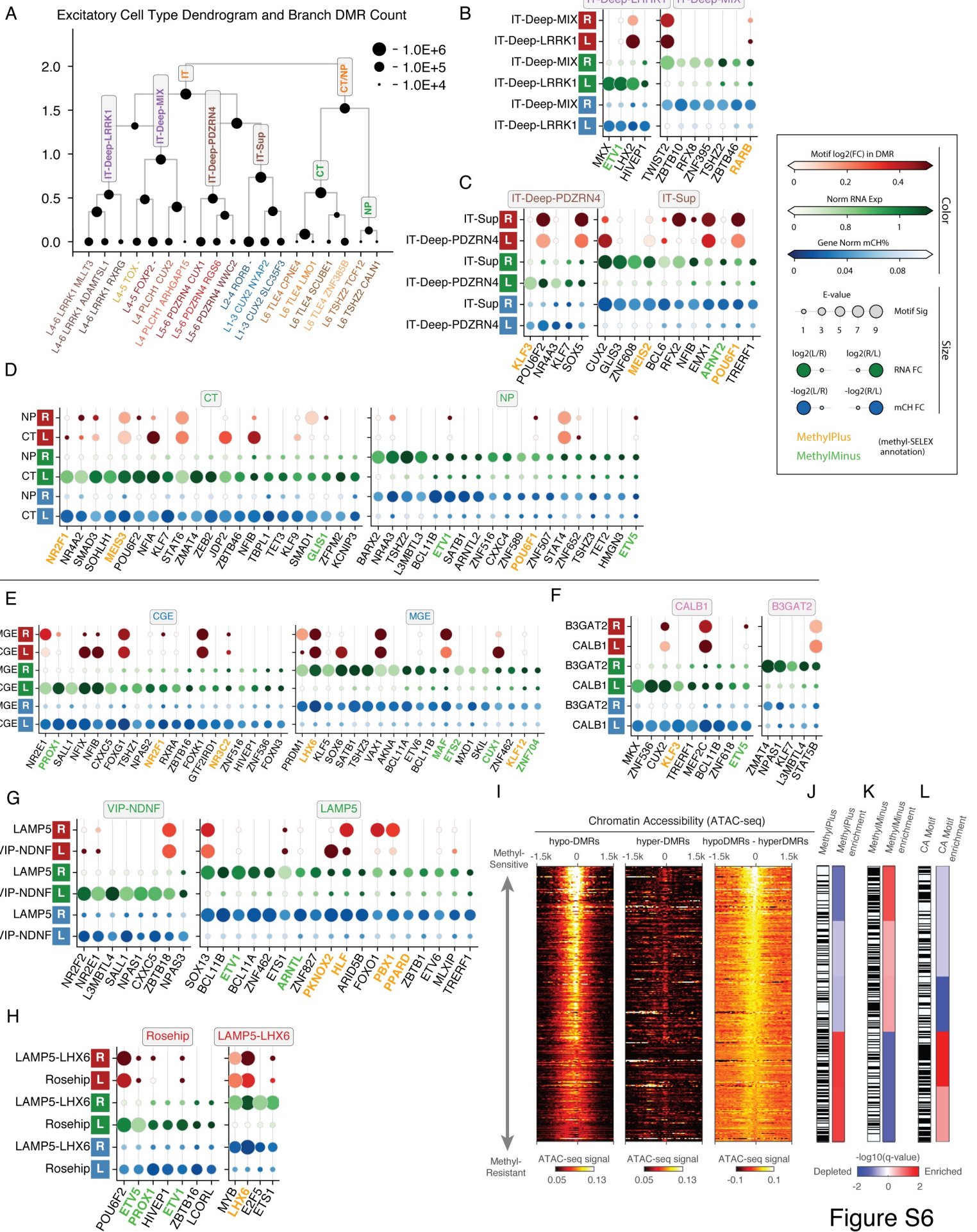


Figure S6

**Figure S6. TF binding motif enrichment across the human cortical neuronal hierarchy, related to Figure 6.** (A) Excitatory neuron subtype dendrogram. The node size represents the number of DMRs detected between the left and right branches. (B-D) Dot plots view for TFs showing lineage-specific motif enrichment, expression and gene body mCH between lineages: intratelencephalic (IT)-Deep-LRRK1 vs IT-Deep-MIX (B); IT-Deep-PDZRN4 vs IT-Sup (C); Cortico-Thalamic projection (CT) vs. layer 5/6 near-projecting neurons (NP) (D). Colors for every two rows from top to top: lineage mean gene body mCH level, lineage mean expression  $\log(1 + \text{CPM})$ , TF motif enrichment  $\log_2(\text{fold change})$ . Sizes for every two rows from bottom to top: relative fold change of mCH level from this branch to the other, relative fold change of expression level, E-value of the motif enrichment test. Colors for the motif names: TF motif methylation preference annotated by methyl-SELEX experiment (Yin et al., 2017), orange indicate MethylPlus, green indicate MethylMinus. (E-H) Dot plot view for TFs showing lineage-specific motif enrichment, expression and gene body mCH between CGE vs. MGE (E), CALB1 vs. B3GAT2 (F), VIP/NDNF vs. LAMP5 (G) and Rosehip vs LAMP5-LHX6 (H). (I) The binding of TFs to hypermethylated regions validated by chromatin accessibility measurement using the snATAC-seq profile. (J-L) Enrichment or depletion of MethylPlus TFs (J), MethylMinus TFs (K) and TFs whose binding motif containing CA dinucleotides (L).



Figure S7



**Figure S7. Prediction of causal cell types for neuropsychiatric traits using partitioned heritability analysis, related to Figure 7.** (A) Partitioned heritability analysis individually comparing (individual models) adult brain cell-type-specific DMRs, fetal cortex lowly methylated regions, and non-brain tissues to the baseline (B) Partitioned heritability analysis (individual models) using adult brain RNA signature genes, ATAC-seq peaks, and NOME-seq peaks. (C-D) Multiple regression partitioned heritability analysis using NOME-seq peaks (C) and gene expression signatures (D). (E-F) Prioritization of neuropsychiatric traits-associated brain cell types using RolyPoly. (G-J) CG methylation and chromatin accessibility profiles at adult brain regulatory elements [DMR (ATAC-pos)] and vestigial enhancers. (K-N) Multiple regression partitioned heritability analysis comparing adult regulatory elements and [DMR (ATAC-pos)] and vestigial enhancers.