

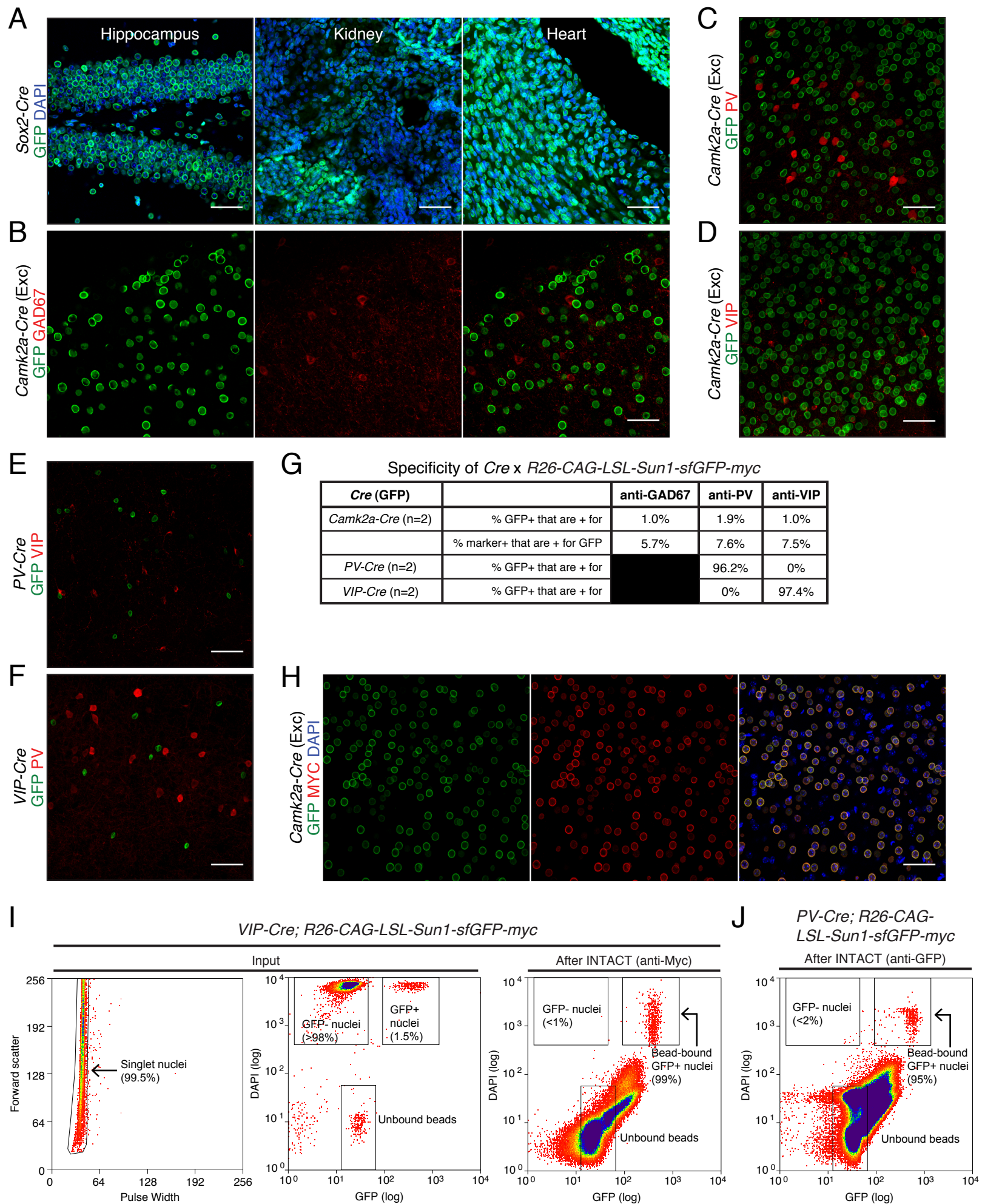
**Neuron**

**Supplemental Information**

## **Epigenomic Signatures**

### **of Neuronal Diversity in the Mammalian Brain**

**Alisa Mo, Eran A. Mukamel, Fred P. Davis, Chongyuan Luo, Gilbert L. Henry, Serge Picard, Mark A. Urich, Joseph R. Nery, Terrence J. Sejnowski, Ryan Lister, Sean R. Eddy, Joseph R. Ecker, and Jeremy Nathans**



### Figure S1. Nuclear Labeling and Specificity of Mouse INTACT, Related to Figure 1

(A) Hippocampus, kidney, and heart from adult *Sox2-Cre; R26-CAG-LSL-Sun1-sfGFP-myc* mice, where Cre recombination occurs at the early embryo stage. Scale bar: 50  $\mu\text{m}$ .

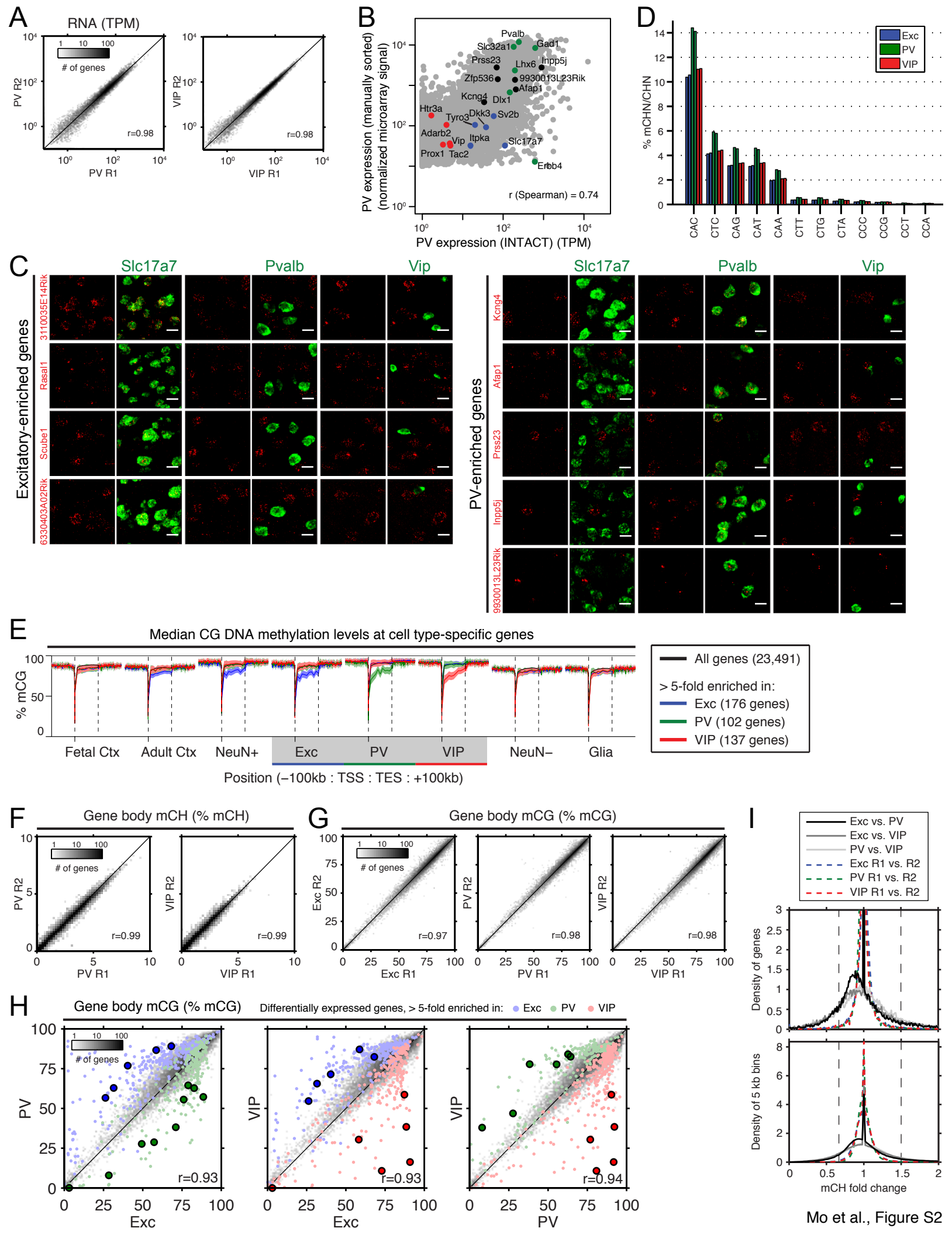
(B) Immunohistochemistry showing that neocortical GFP<sup>+</sup> nuclei in *Camk2a-Cre; R26-CAG-LSL-Sun1-sfGFP-myc* adult mice do not co-localize with GAD67, an inhibitory neuron marker. Scale bar: 50  $\mu\text{m}$ .

(C-F) Immunohistochemistry showing that neocortical GFP<sup>+</sup> nuclei in *Camk2a-Cre; R26-CAG-LSL-Sun1-sfGFP-myc* adult mice do not co-localize with PV (C) or VIP (D). Similarly, neocortical GFP<sup>+</sup> nuclei in *PV-Cre; R26-CAG-LSL-Sun1-sfGFP-myc* adult mice do not co-localize with VIP (E), and neocortical GFP<sup>+</sup> nuclei in *VIP-Cre; R26-CAG-LSL-Sun1-sfGFP-myc* adult mice do not co-localize with PV (F). Scale bars: 50  $\mu\text{m}$ .

(G) Quantification of Cre driver specificity by immunostaining. Each Cre driver was crossed with *R26-CAG-LSL-Sun1-sfGFP-myc* mice, and the percent of GFP<sup>+</sup> cells that co-localize with the indicated markers was quantified. For the *Camk2a-Cre* driver, the percentage of GAD67, PV, and VIP cells that co-localize with GFP was also quantified; furthermore, quantification of GFP and NeuN staining (Figure 1B) showed that 100% of GFP<sup>+</sup> nuclei were also NeuN<sup>+</sup>. Counts were made using 100  $\mu\text{m}$  vibratome sections and >200 nuclei per mouse (n=2).

(H) Myc labeling co-localizes with GFP labeling in the neocortex of adult *Camk2a-Cre; R26-CAG-LSL-Sun1-sfGFP-myc* mice. Scale bar: 50  $\mu\text{m}$ .

(I and J) Assessment of INTACT purification by flow cytometry. INTACT was performed using anti-Myc to isolate nuclei from the neocortices of two *VIP-Cre; R26-CAG-LSL-Sun1-sfGFP-myc* mice (I). Analysis of input nuclei (after step 2 in Figure 1C) shows that 99.5% of input nuclei are singlets (left), and 1.5% of input nuclei are GFP<sup>+</sup> (middle). Unbound beads remaining after the pre-clear step were identified using a beads-only control (data not shown). After INTACT purification, 99% of bead-bound nuclei are GFP<sup>+</sup> (right). Because multiple magnetic beads are bound to each GFP<sup>+</sup>/Myc<sup>+</sup> nucleus, the DAPI fluorescence is variably reduced relative to input nuclei. INTACT was performed using anti-GFP to isolate nuclei from the neocortices of two *PV-Cre; R26-CAG-LSL-Sun1-sfGFP-myc* mice (J). After INTACT purification, 95% of bead-bound nuclei are GFP<sup>+</sup>. The percentages of singlet nuclei, GFP<sup>-</sup> nuclei, and GFP<sup>+</sup> nuclei were determined by the gates outlined in black in each plot.



## Figure S2. Gene Expression and DNA Methylation Analysis, Related to Figure 2

(A) Pairwise comparisons of gene expression levels between replicates in PV (left) and VIP (right) neurons.  $r$ , Pearson correlation of  $\log(\text{TPM}+0.1)$ . TPM, transcripts per million.

(B) Scatterplot showing high correlation between gene expression of INTACT-purified (*PV-Cre; R26-CAG-LSL-Sun1-sfGFP-myc*) and manually-sorted (G42 transgenic; Okaty et al., 2009) PV neurons. Selected cell type-specific genes are labeled (blue, Exc; green, PV; red, VIP) as well as candidate PV-enriched genes (black) tested by *in situ* hybridization.  $r$ , Spearman correlation.

(C) Double fluorescent ISH showing correct co-localization for nine genes predicted to be enriched in excitatory (left) or PV (right) neurons. *Slc17a7*, *Pvalb*, and *Vip* mark excitatory, PV, and VIP neurons, respectively. A 10th probe (*Zfp536*) did not co-localize with *Slc17a7*, *Pvalb*, or *Vip* at our level of detection (data not shown), and probe labeling was presumably in oligodendrocytes (Dugas et al., 2006).

(D) Barplot showing % mC for each non-CG methylation trinucleotide context.

(E) Median  $\pm$  1 SEM of % mCG within and surrounding gene bodies, showing an inverse correlation between expression and DNA methylation at differentially expressed genes determined from our RNA-seq data (>5 fold-change for one cell type relative to both of the other cell types). TSS, transcription start site; TES, transcription end site; SEM, standard error of the mean.

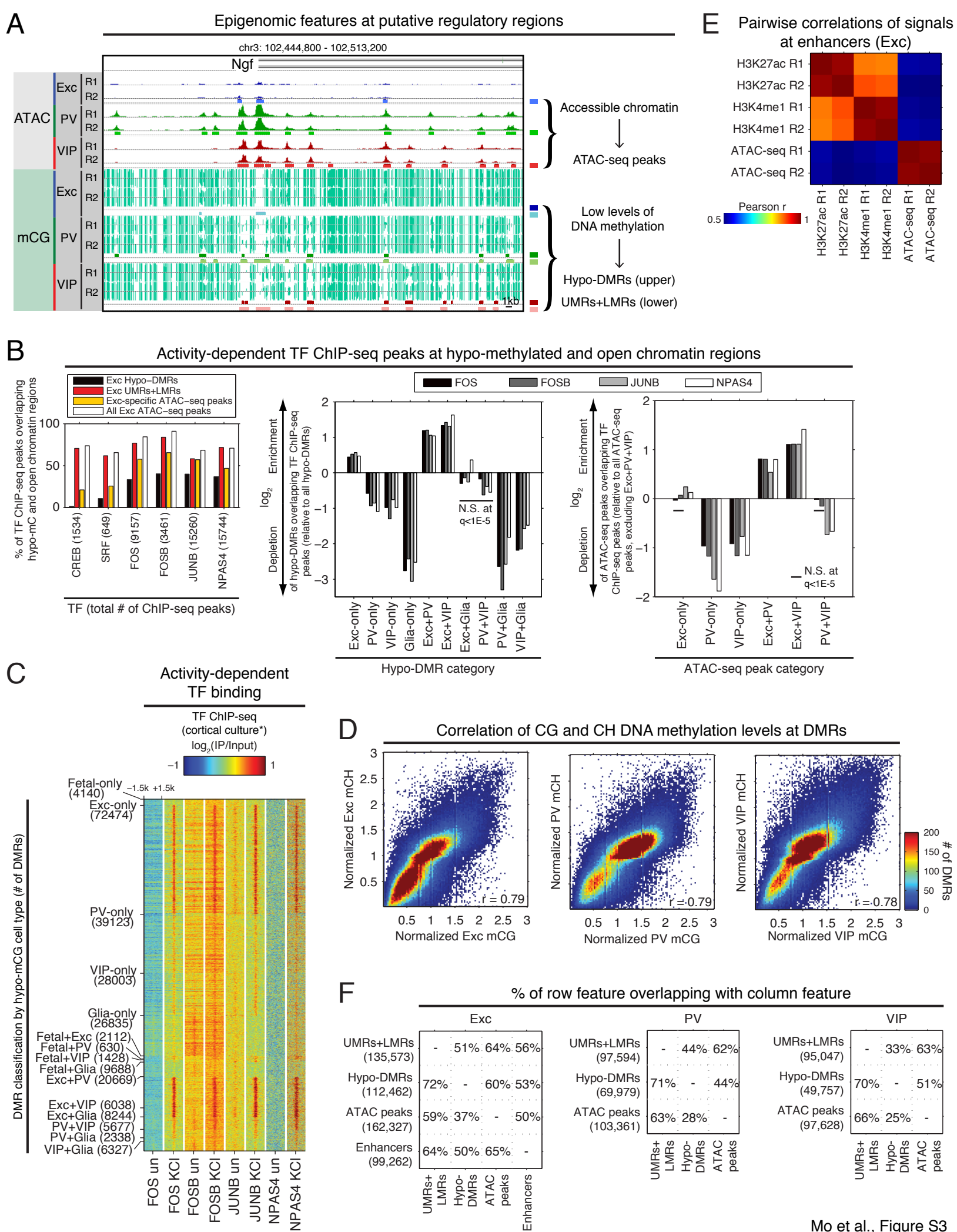
(F) Pairwise comparisons of gene body % mCH between replicates in PV (left) and VIP (right) neurons.

(G) Pairwise comparisons of gene body % mCG between replicates in excitatory (left), PV (middle), and VIP (right).

(H) Pairwise comparisons of gene body % mCG across cell types. Colored dots correspond to the same genes shown in Figure 2B.

(I) Density plots showing ratios of CH methylation in gene bodies (top) and in 5 kb genomic bins (bottom) across cell types and between replicates. Each distribution was normalized by the median ratio. Dotted lines are at 0.67 and 1.5.





**Figure S3. Correlations Across Epigenomic Marks and Relevance of Neuron Subtype-Specific Hypo-Methylation and Open Chromatin to Induced Neuronal Activity, Related to Figure 3**

(A) Examples of regulatory elements marked by accessible chromatin (peaks in ATAC-seq read density, upper tracks) and low levels of DNA methylation (hypo-DMRs and UMRs+LMRs, lower tracks) near *Ngf*. Locations of ATAC-seq peaks, hypo-DMRs, and UMRs+LMRs are shown below the corresponding raw reads. R1, replicate 1; R2, replicate 2.

(B) Barplot showing that the majority of binding sites of six activity-dependent TFs in KCl-depolarized cortical cultures (Kim et al., 2010; Malik et al., 2014) overlap with excitatory neuron UMRs+LMRs and all excitatory ATAC-seq peaks (left). Binding sites for FOS, FOSB, JUNB, and NPAS4 also overlap extensively with excitatory-specific hypo-DMRs and ATAC-seq peaks. The number of total ChIP-seq peaks for each TF is shown in parentheses. Barplot showing the enrichment and depletion of each hypo-DMR category overlapping TF ChIP-seq peaks (middle). CREB and SRF were excluded since their enrichments and depletions were insignificant at  $q < 1E-5$ . Barplot showing the enrichment and depletion of each differential ATAC-seq peak category overlapping TF ChIP-seq peaks (right).

(C) At the same regions as in Figure 3E (i.e., 3 kb windows centered at DMRs hypo-methylated in one or two cell types), heatmap showing TF ChIP-seq reads from unstimulated (un) and depolarized (KCl) cortical cultures (TF ChIP-seq from Kim et al., 2010; Malik et al., 2014).

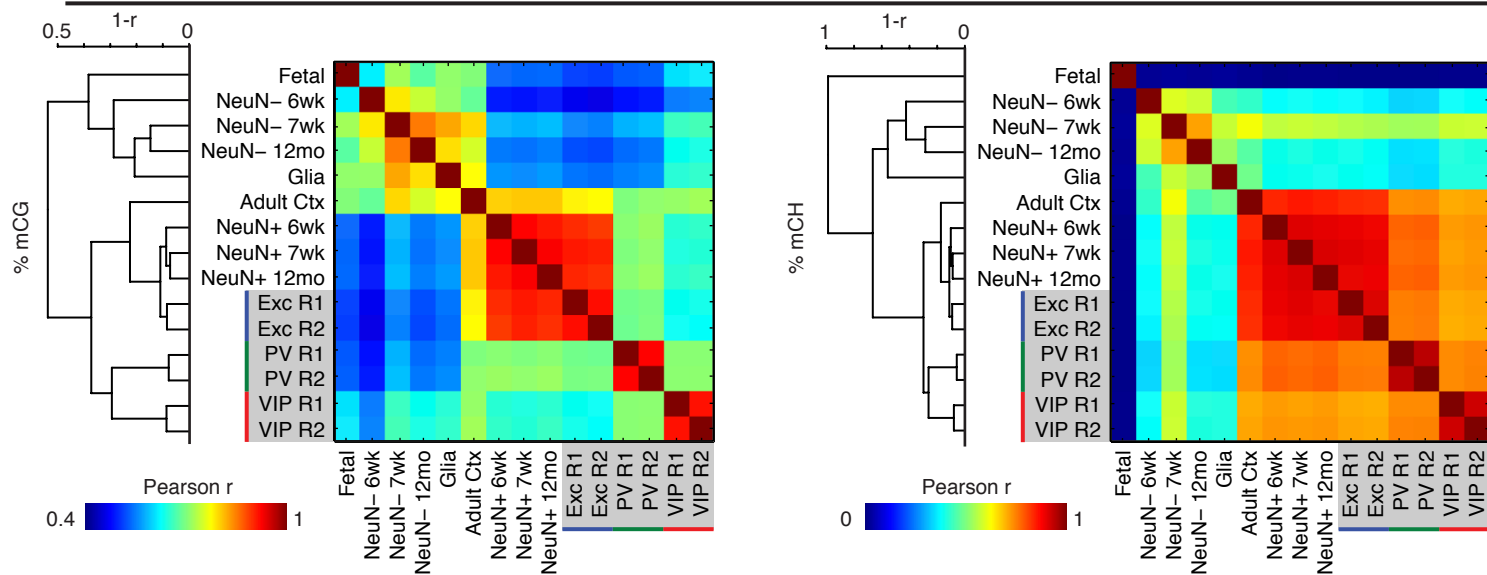
(D) Scatterplots showing high correlation between mCG and mCH at DMRs. mCG and mCH levels in each DMR were normalized by the mean mCG and mCH in that DMR across the three cell types.  $r$ , Pearson correlation.

(E) A matrix showing pairwise Pearson correlations of H3K27ac, H3K4me1, and sub-nucleosomal (<100 bp) ATAC-seq reads at putative enhancers. H3K4me1 and H3K27ac signals are generally well-correlated at a global level; however, individual enhancers can be poised (H3K4me1+; 42,540 enhancers) or active (H3K4me1+/H3K27ac+; 48,781 enhancers). ATAC-seq signal is also correlated, albeit to a lesser degree ( $r \sim 0.5$ ), with both H3K4me1 and H3K27ac signal at enhancers. R1, replicate 1; R2, replicate 2.

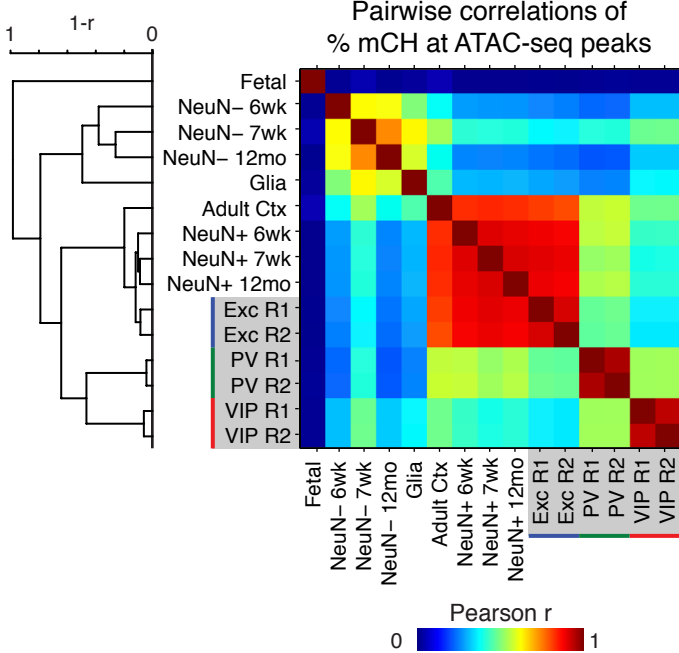
(F) In each cell type, the majority of hypo-DMRs are a subset of UMRs+LMRs. The majority of ATAC-seq peaks and UMRs+LMRs overlap. For excitatory neurons, approximately half of ATAC-seq peaks, hypo-DMRs, and UMRs+LMRs overlap with enhancers identified using histone modifications. Although these DNA methylation and chromatin features are overlapping, they are not synonymous. Part of the difference may arise from statistical thresholds set in the identification of each region; however, each type of dataset also provides non-redundant and complementary information that depend on the genomic context. The numbers in parentheses indicate the total number of these features identified in each cell type.

**A**

## Pairwise correlations of % mCG and % mCH in 500 bp genomic bins

**B**

## Pairwise correlations of % mCH at ATAC-seq peaks

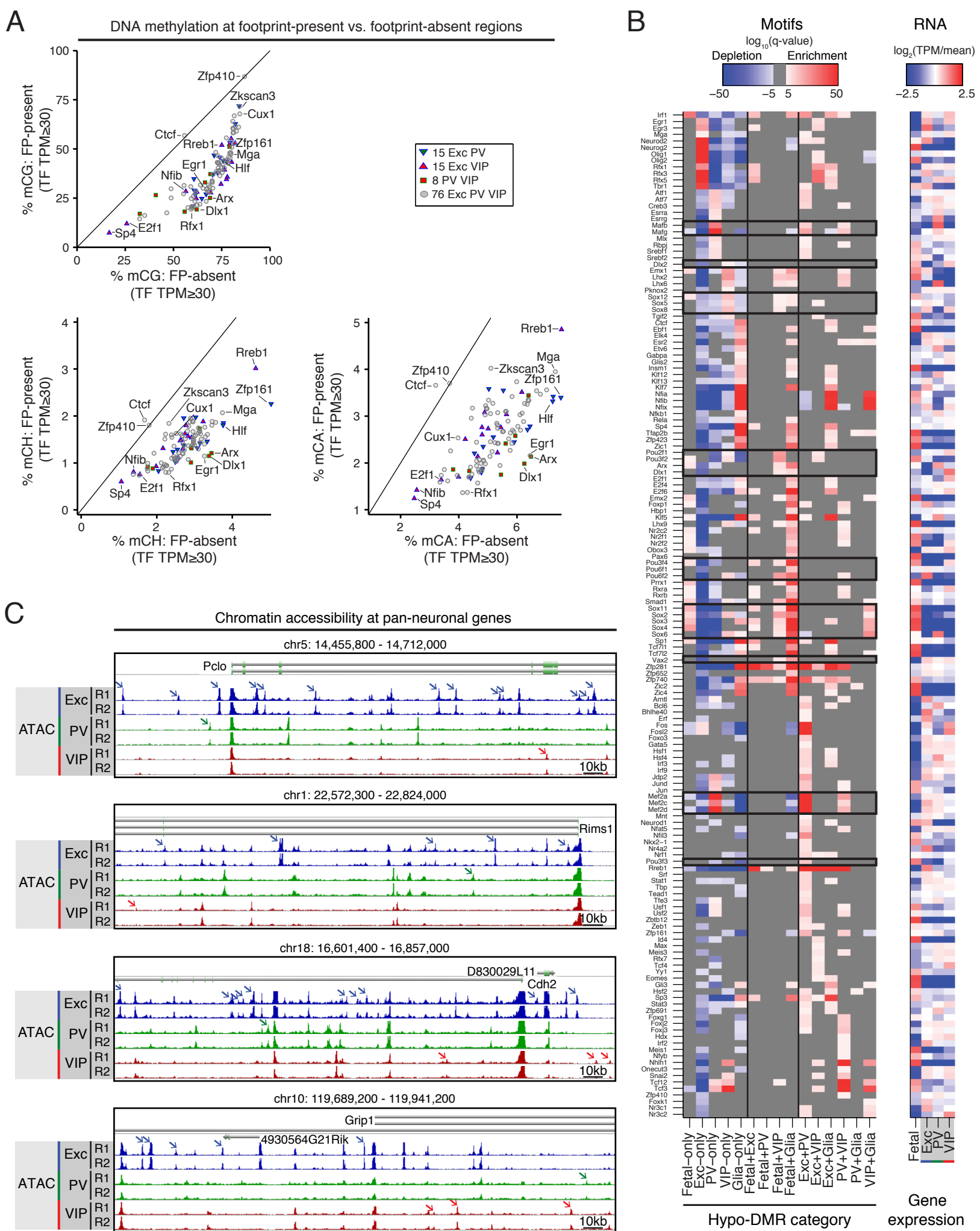




**Figure S4. Epigenomic Correlations Across Cell Types and Development, Related to Figure 4**

(A) Matrices showing pairwise Pearson correlations for % mCG (left) and for % mCH (right) in 500 bp genomic bins across all autosomes. Dendrograms show hierarchical clustering using complete linkage and 1-Pearson correlation as the metric.

(B) A matrix showing pairwise Pearson correlations for % mCH at ATAC-seq peaks. The dendrogram shows hierarchical clustering using complete linkage and 1-Pearson correlation as the metric.

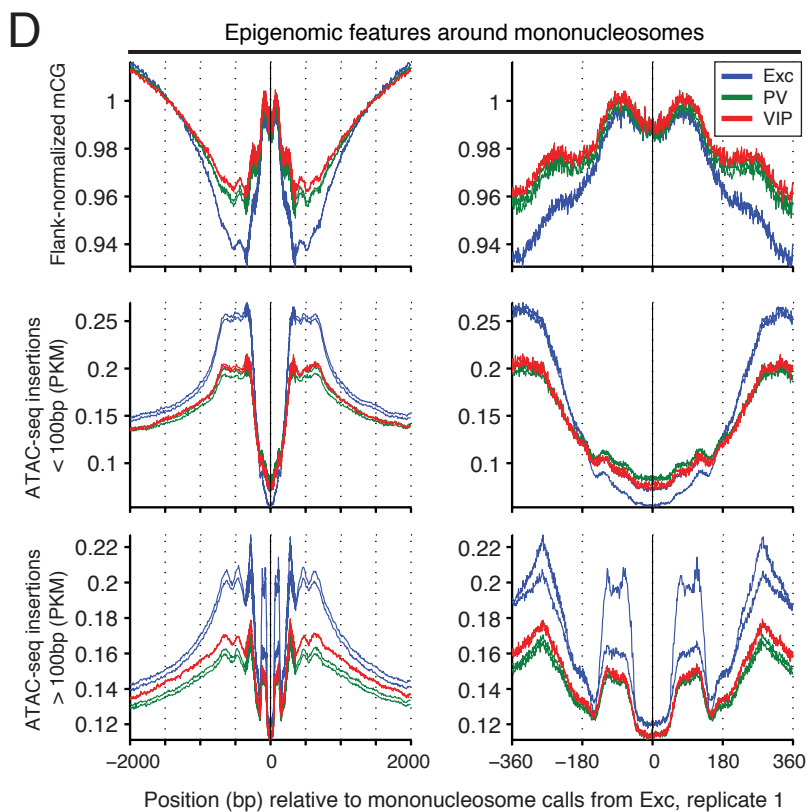
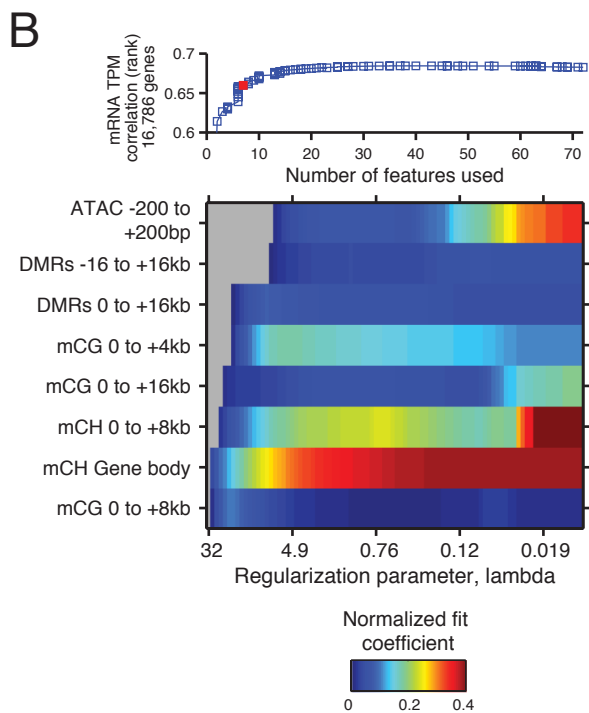
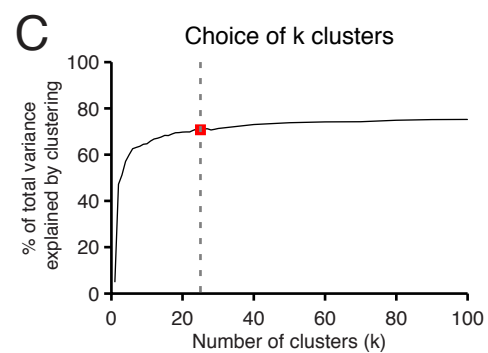
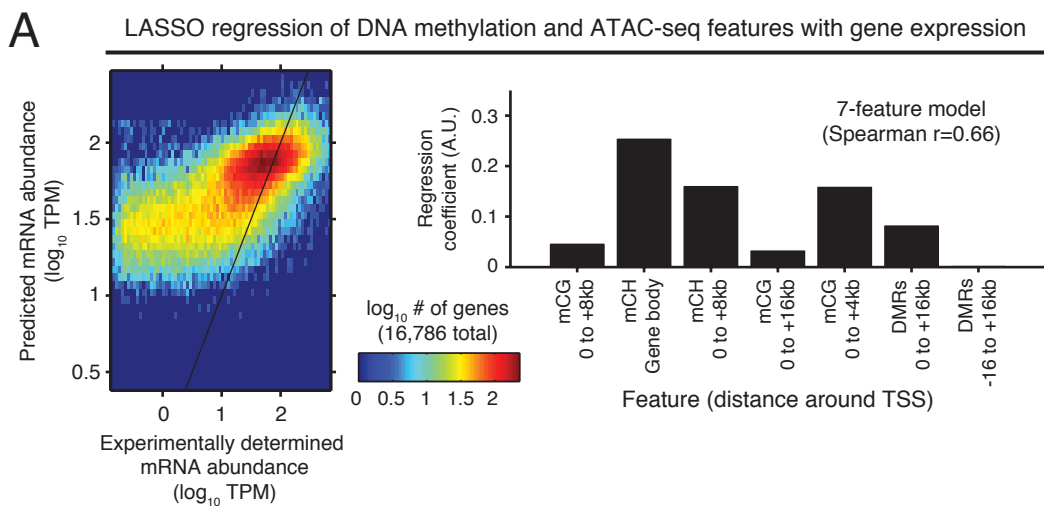


**Figure S5. Analysis of Putative TF Binding at Neuronal Regulatory Regions, Related to Figure 5**

(A) Scatterplots showing, for expressed TFs ( $\text{TPM} \geq 30$ ), % mCG (top), % mCH (bottom left), and % mCA (bottom right) around regions that are footprinted in one cell type (y-axis) versus regions that are not footprinted in that cell type, but are footprinted in a different cell type (x-axis). Most TF footprints lie in regions of lower DNA methylation, relative to the methylation levels found in cell types without footprints for the same regions. Exceptions include CTCF and ZFP410.

(B) Heatmap showing TF motif enrichments (left) and gene expression (right) for all categories of DMRs that are hypo-methylated in one or two cell types across excitatory, PV, and VIP neurons as well as glia and fetal cortex. Boxes indicate TFs mentioned in the main text.

(C) Examples of pan-neuronal genes (from Hobert et al., 2010) surrounded by cell type-specific and pan-neuronal regions of increased chromatin accessibility, as determined by peaks of ATAC-seq read density. Arrows point to a subset of cell type-specific ATAC-seq peaks. R1, replicate 1; R2, replicate 2.



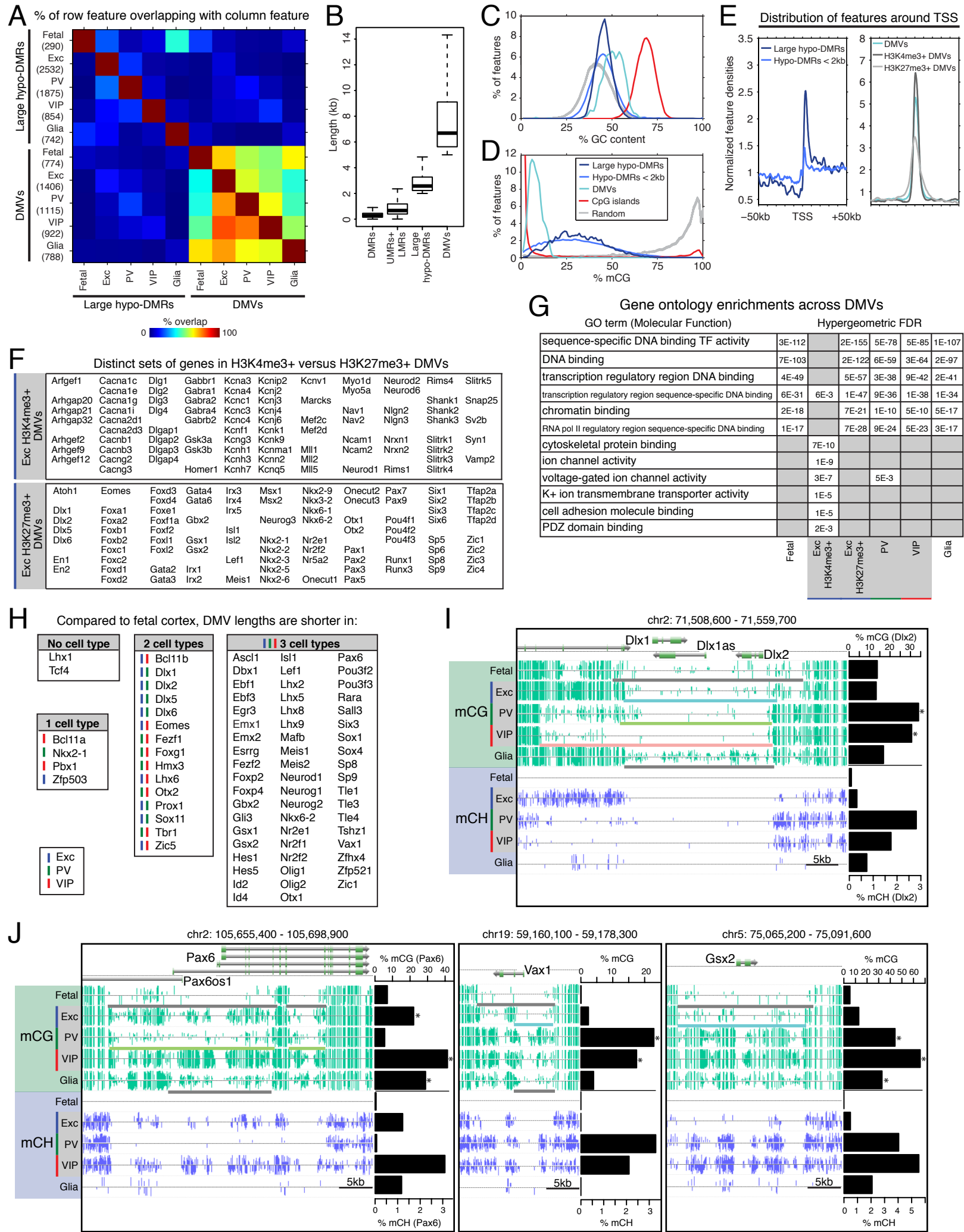
**Figure S6. Integrative Analysis of Epigenomic Features, Related to Figure 6**

(A) LASSO regression using the top 7 selected epigenomic features gives a Spearman correlation of 0.66, with intragenic non-CG methylation as the most informative feature. Epigenomic features used in the regression were mCG, mCH, ATAC-seq, and DMR density at different positions around genes (see Table S5). A.U., arbitrary units.

(B) Line plot showing that LASSO regression using more than ~7 features does not generate substantially higher correlations (top). The normalized fit coefficient for the 8 best features is shown as a function of the regularization parameter (bottom). The red square indicates 7 features.

(C) Choice of the number of clusters used for k-means clustering.

(D) Line plots showing lower mCG and ATAC-seq read density at the mononucleosome core.





## Figure S7. Large Hypo-Methylated Domains, Related to Figure 7

(A) Large hypo-DMRs are generally non-overlapping across cell types, whereas DMVs show high overlap across cell types. Large hypo-DMRs and DMVs are generally non-overlapping regions in the same cell type. The numbers of large hypo-DMRs and DMVs identified in each cell type are indicated in parentheses.

(B) Boxplot showing the length distributions for large hypo-methylation features compared to all DMRs and UMRs+LMRs. Large hypo-DMRs and DMVs are both multi-kilobase DNA methylation features. By definition, the lower size limits are 2 kb for large hypo-DMRs and 5 kb for DMVs. Autosomal features for excitatory, PV, and VIP neurons were combined. Outliers are omitted in the graphical representation.

(C-D) Distribution of GC content (C) and CG methylation level (D) across DNA methylation features. Excitatory neuron features and methylation levels were used, as well as randomly selected genomic regions matching the sizes of excitatory hypo-DMRs with lengths less than 2 kb.

(E) Line plots showing that large hypo-DMRs are enriched downstream of the TSS whereas DMVs are enriched equally across the TSS. Excitatory neuron features were used.

(F) Representative selection of genes in excitatory DMVs that overlap H3K4me3+ peaks (top) and H3K27me3+ domains (bottom).

(G) Gene ontology (GO) categories (McLean et al., 2010) related to transcription regulation and TF activity are strongly enriched at H3K27me3+ excitatory DMVs and DMVs in other cell types, including fetal brain. H3K4me3+ excitatory DMVs are enriched for terms related to mature neuronal function.

(H) Out of 77 developmental TFs (Visel et al., 2013) that overlap fetal DMVs, the DMV lengths for 75 TFs are shorter in at least one adult cell type relative to fetal cortex. For each TF, the cell type(s) with decreased DMV length(s) are indicated. See Table S6.

(I-J) (I) DNA methylation levels for a region around *Dlx1/2*, showing extensive neuron subtype-specific differences in the boundaries of DMVs that correlate with developmental shifts in the expression of *Dlx2* and *Dlx1*. (J) DNA methylation levels for a region around *Pax6* (left), *Vax1* (middle), and *Gsx2* (right). *Pax6* is expressed during excitatory neuron development and in the caudal ganglionic eminence (birthplace of VIP neurons), whereas *Vax1* and *Gsx2* are expressed during inhibitory neuron development. Expression levels of all three TFs are largely down-regulated in mature neurons. For (I) and (J), see Table S6 for annotations and references. Barplots show the % mCG and % mCH for each cell type at the region between the dotted lines in Figure 7E. \*  $q < 1 \times 10^{-10}$  (mCG, adult cell type compared to fetal cortex, 1-sided FET with Benjamini-Hochberg correction). In the browser representation, light-colored bars indicate DMVs.

**Table S1. Sample and Sequencing Information, Related to Figure 1**

Sample and sequencing information for each dataset

**Table S2. Gene Expression Levels and Identifications of Differentially Expressed Genes, Related to Figure 2**

Pairwise identifications of differentially expressed genes and full expression table

**Table S3. Lists of DMRs, UMRs+LMRs, and ATAC-Seq Peaks, Related to Figure 3**

Identification and classification of DMRs, UMRs+LMRs, and ATAC-seq peaks

**Table S4. Lists of TFs Predicted to Regulate Cell Type-Specific Gene Expression, Plus TF mC Sensitivity and TF-TF Regulatory Networks, Related to Figure 5**

TFs that are predicted to regulate cell type-specific gene expression; mC sensitivity of TFs; TF-TF regulatory networks

**Table S5. Coordinates of Windows Used for Generalized Linear Model (in Kb) and GO Analysis for K-Means Clusters, Related to Figure 6**

Window coordinates around TSS, gene body, and gene body flanking regions; Enrichment of GO terms for k-means clusters

**Table S6. Lists of Large Hypo-DMRs and DMVs, Plus Analysis of DNA Methylation Dynamics at DMVs, Related to Figure 7**

For each cell type, coordinates for large hypo-DMRs and DMVs as well as analysis of developmental DMV changes

## Supplemental Experimental Procedures

### Generation of the *R26-CAG-LSL-Sun1-sfGFP-Myc* mouse line

Using the approach of Henry et al., 2012, we tagged the C-terminus of mouse SUN1 by attaching two copies of superfolder GFP, a variant of GFP with increased brightness and stability (Pédélec et al., 2006), and six tandem copies of Myc. We inserted this cassette into a *Rosa26* targeting vector (Soriano, 1999) downstream of a *CAG* promoter and a *loxP-3x polyA-loxP* sequence. The construct was electroporated into 129-derived ES cells, and correctly targeted cells were injected into C57BL/6J blastocysts to screen for chimeras. Chimeric males were bred to C57BL/6J females and intercrossed to obtain homozygotes. *R26-CAG-LSL-Sun1-sfGFP-Myc* mice have been deposited at JAX (Stock 021039).

### Mouse lines

*Camk2a-Cre* (Stock 005359), *PV-Cre* (Stock 008069), *VIP-Cre* (Stock 010908), and *Sox2-Cre* (Stock 008454) mice were obtained from JAX.

### Immunohistochemistry and microscopy

Mice were anesthetized with ketamine/xylazine, perfused with 4% paraformaldehyde (PFA), and post-fixed for 1 hour at room temperature. Brains were sectioned at 100  $\mu\text{m}$  using a vibratome. Sections were blocked with 10% NGS and 0.25% Triton in PBS and incubated with the following antibodies overnight at 4°C: rabbit anti-Parvalbumin (1:5000, Swant PV 25), rabbit anti-Vasoactive intestinal peptide (1:200, ImmunoStar 20077), mouse anti-NeuN (1:500, Millipore MAB377). Either chicken anti-GFP (1:500, Aves GFP-1020), rabbit anti-GFP (1:400, Life Technologies A11122), or rabbit anti-Myc (1:50,000, homemade) was co-incubated. For staining with anti-VIP, mice were perfused with 2% PFA as heavy fixation decreased the VIP signal. For mouse anti-GAD67 (1:800, Millipore MAB5406), no Triton was included in any step, and both primary and secondary antibody incubations were performed at room temperature for 36 hours. For fluorescent labeling, the sections were incubated with Alexa Fluor-conjugated IgG secondary antibodies (1:400, Life Technologies) and DAPI before mounting with Fluoromount G (SouthernBiotech). For assessment of Cre driver specificity, we counted more than 200 neocortical nuclei for each mouse and two mice per Cre driver. Images were taken using a Zeiss LSM700 confocal microscope (immunohistochemistry, fluorescent *in situ* hybridization) or a Zeiss Imager Z1 and Apotome system (bead-bound nucleus). Image processing was performed using ImageJ and Adobe Photoshop, including adjustments of brightness, contrast, and levels in individual color channels for merged images.

### Mouse INTACT procedure

For each experiment, the neocortices of one to two mice were rapidly dissected in ice-cold homogenization buffer (0.25M sucrose, 25mM KCl, 5mM MgCl<sub>2</sub>, 20mM Tricine-KOH). The tissue was minced with a razor blade and Dounce homogenized using a loose pestle in 5 mL of homogenization buffer supplemented with 1mM DTT, 0.15mM spermine, 0.5mM spermidine, and EDTA-free protease inhibitor (Roche 11 836 170 001). A 5% IGEPAL-630 solution was added to bring the homogenate to 0.3% IGEPAL-630, and the homogenate was further dounced with five strokes of the tight pestle. When purifying RNA, RNasin Plus RNase Inhibitor (Promega N2611) was added at 60 U/mL. The sample was filtered through a 40  $\mu\text{m}$  strainer (Fisher Scientific 08-771-1), mixed with 5 mL of 50% iodixanol density medium (Sigma D1556), underlayered with a gradient of 30% and 40% iodixanol, and centrifuged at 10,000g for 18 minutes in a swinging bucket centrifuge at 4°C. Nuclei were collected at the 30%-40% interface and pre-cleared by incubating with 20  $\mu\text{L}$  of Protein G Dynabeads (Life Technologies 10003D) for 10 minutes. After removing the beads with a magnet, the mixture was diluted with wash buffer (homogenization buffer plus 0.4% IGEPAL-630) and incubated with 10  $\mu\text{L}$  of 0.2 mg/mL rabbit

monoclonal anti-GFP antibody (Life Technologies G10362) or anti-Myc antibody (homemade, ~2 µg) for 30 minutes. 60 µL of Dynabeads were added, and the mixture was incubated for an additional 20 minutes. To increase yield, the bead-nuclei mixture was placed on a magnet for 30 seconds to 1 minute, completely resuspended by inversion, and placed back on the magnet. This was repeated 5-7 times. Bead-bound nuclei were passed through a 20 µm strainer (Partec 04-0042-2315) and washed with 2 x 10 mL, 1 x 2.5 mL, and 1 x 1 mL wash buffer. All steps were performed on ice or in the cold room, and all incubations were carried out using an end-to-end rotator.

All calculations of INTACT specificity and yield used pooled neocortices (approximately dorsal 2/3 of cortex) of two 8-11 week old mice. To calculate the specificity of mouse INTACT, bead-bound nuclei were stained with DAPI, viewed by fluorescence microscopy, and the numbers of GFP+ and GFP- nuclei were counted (100-200 nuclei per experiment). To calculate the yield of mouse INTACT, input nuclei (i.e., after step 2 in Figure 1C) and bead-bound nuclei were stained with DAPI. The yield was determined from the total number of input nuclei, the % of GFP+ nuclei in the input, and the total number of bead-bound nuclei after INTACT purification (all quantified by fluorescence microscopy or hemocytometer, 100-200 nuclei per experiment).

### **Flow cytometry**

Beads-only control, input nuclei, and bead-bound INTACT-purified nuclei (using anti-Myc antibody) from *VIP-Cre; R26-CAG-LSL-Sun1-sfGFP-Myc* neocortices as well as beads-only control, input nuclei, and bead-bound INTACT-purified nuclei (using anti-GFP antibody) from *PV-Cre; R26-CAG-LSL-Sun1-sfGFP-Myc* neocortices were analyzed using a MoFlo MLS high-speed cell sorter (Beckman Coulter).

### **Extraction of RNA, DNA, and native nucleosomes**

Bead-bound nuclei or whole neocortical nuclei were directly resuspended in Buffer RLT for RNA purification using the RNeasy Micro kit (Qiagen 74004) following the standard protocol with on-column DNase digestion. For RNA preparation from whole neocortical nuclei, nuclei were prepared identically to the INTACT procedure, except that the 40% iodixanol layer was omitted, and nuclei were pelleted by centrifugation and resuspended in Buffer RLT. Bead-bound nuclei were resuspended in PBS for DNA purification using the DNeasy Blood and Tissue kit (Qiagen 69504). Nucleosomes for native ChIP-seq were prepared as previously described (Henry et al., 2012). Briefly, 1-2 million bead-bound nuclei were digested with 0.025 units/µL micrococcal nuclease (Worthington LS004798) in 500 µL of 15mM HEPES pH 7, 1mM KCl, 2mM MgCl<sub>2</sub>, 2mM CaCl<sub>2</sub>, 340mM sucrose, 0.15mM spermine, 0.5mM spermidine, and 5mM sodium butyrate at 37°C for 15 minutes. The reaction was terminated by the addition of EGTA to 2mM final concentration. Nucleosomes were extracted for 30 minutes on ice with 200 µL 15mM HEPES pH7, 200mM NaCl, 25mM KCl, 2mM MgCl<sub>2</sub>, 1mM EGTA, 340mM sucrose, 0.15mM spermidine, 0.15mM spermine, and 5mM sodium butyrate. A second 30 minute extraction was performed with the same buffer except the salt concentration was raised to 400mM NaCl. The extracts were combined and dialyzed overnight against 15mM HEPES pH7, 25mM KCl, 1mM β-mercaptoethanol, 1mM PMSF, and 5mM sodium butyrate using a 10K cut-off Slide-a-Lyzer dialysis device (Thermo Scientific 88401).

### **RNA-seq library construction and sequencing**

RNA quality was measured by an Agilent Bioanalyzer, with RIN scores consistently greater than 8. Total RNA (2-50 ng) was converted to cDNA and amplified using Nugen Ovation RNA-seq System V2 (Nugen 7102). No selection for poly-adenylated RNA was used. All RNA samples received a 1:10,000 dilution of ERCC RNA (Life Technologies 4456740). Amplified cDNA was fragmented, end-repaired, linker adapted, and sequenced for 50 cycles on an Illumina HiSeq 2500 instrument. Image analysis and base calling were performed with the standard Illumina pipeline versions RTA 1.12.4.2 and 1.17.20.

## **MethylC-seq library construction and sequencing**

MethylC-seq libraries were constructed as previously described (Lister et al., 2013), except that samples were PCR amplified with KAPA HiFi HotStart Uracil+ ReadyMix (Kapa Biosystems KK2802) using the following PCR conditions: 2 minutes at 95°C, 4 cycles of [15 seconds at 98°C, 30 seconds at 60°C, 4 minutes at 72°C], and 10 minutes at 72°C. Libraries were sequenced on an Illumina HiSeq 2000 up to 101 cycles. Image analysis and base calling were performed with the standard Illumina pipeline version RTA 2.8.0.

## **ATAC-seq library construction and sequencing**

Approximately 50,000 bead-bound nuclei were transposed in a 50 µL volume of 1X TD buffer and 2.5 µL Tn5 transposase (Illumina FC-121-1030) for 30 minutes at 37°C, as previously described (Buenrostro et al., 2013), with the modification that fragmented genomic DNA was recovered using Buffer QG coupled with MinElute spin columns (Qiagen 28604). Transposed genomic DNA was amplified by five cycles of quantitative PCR. 10% of the PCR was subjected to an additional 20 cycles of SYBR green-based qPCR while the remainder of the sample was left on ice. Analysis of the qPCR data allowed a rough estimate of the number of additional cycles needed to generate product at 25% saturation. Typically, four to seven additional PCR cycles were added to the initial set of five cycles. Amplified DNA was purified on AMPure XP beads (Beckman A63881), analyzed on an Agilent Bioanalyzer, and sequenced (paired-end) on an Illumina HiSeq 2500 for 101 cycles. Image analysis and base calling were performed with the standard Illumina pipeline versions RTA 1.17.20 and 1.17.21.3.

## **ChIP-seq library construction and sequencing**

We used the HT ChIP-seq protocol (Garber et al., 2012) for the ChIP reactions and subsequent library construction with the following modifications. For each reaction, chromatin prepared from 0.5-1 million nuclei was incubated with 1 µg antibody and 25 µL Protein G Dynabeads. The following antibodies were used: rabbit anti-H3K27me3 (Millipore 07-449), rabbit anti-H3K27ac (Abcam ab4729), rabbit anti-H3K4me3 (Abcam ab8580), and rabbit anti-H3K4me1 (Abcam ab8895). ChIP-enriched and input DNA was end-repaired, linker adapted, amplified, and sequenced on an Illumina HiSeq 2500 for 50 cycles. Image analysis and base calling were performed with the standard Illumina pipeline version RTA 1.17.20.

## **Fluorescent *in situ* hybridization**

cDNA libraries were generated from neocortical 8 week old C57BL/6J brains with the SuperScript III First-Strand Synthesis System (Life Technologies 18080-051). The following primers were used for producing probes: *3110035E14Rik* (For: 5'-GATAAGAAAGCACTGTGGTCCC-3', Rev: 5'-ACAGTGAGAAAATCCACCCAAG-3'); *Rasall* (For: 5'-GTGTGTTCTGGGGCAACC, Rev: 5'-GCTTCTCCACACACCGCT-3'); *Scubel1* (For: 5'-TGGACTAGGTGTTGTGTGGAAG-3', Rev: 5'-TAGCTTCTCCCTGAGTTCCAAG-3'); *6330403A02Rik* (For: 5'-GGCATGCTTATCCAACACTACACA-3', Rev: 5'-TACATTTTCATGAGTCCCAGTGC-3'); *Kcng4* (For: 5'-CCATCCCATGGCTGAGAC-3', Rev: 5'-CAGCATTAGCCCCATTG-3'); *Afap1* (For: 5'-CAGCAAGGCACAGACCCT-3', Rev: 5'-TGACTGCTGGGAGCCTTC-3'); *Prss23* (For: 5'-GGGGCAGGATCCACTTCT-3', Rev: 5'-AGCAGCGTGGGAATTCTG-3'); *Inpp5j* (For: 5'-CTTCAACTTTGTGCTGGTGAG-3', Rev: 5'-GTAACCCAGAATGAAGTCTCCG-3'); *9930013L23Rik* (For: 5'-ATCTGGGTGACTCTGGAGAC-3', Rev: 5'-AGAGGCCACCTCTTCTCTC-3'); *Zfp536* (For: 5'-TATCAGGCCTGGCAGCTC-3', Rev: 5'-AGTCGATTCCGGGGAGAC-3'); *Slc17a7* (For: 5'-CAGAGCCGGAGGAGATGA-3' ; Rev: 5'-TTCCCTCAGAAACGCTGG-3'); *Pvalb* (For: 5'-TCTGCTCATCCAAGTTGCAG-3' ; Rev: 5'-TCCTGAAGGACTCAACCCC-3'); *Vip* (For: 5'-CCTTCCCTAGAGCAGAACTTCAG-3' ; Rev: 5'-ACATCAATTTTCTCGATTGCTAC-3'). For all genes except *9930013L23Rik*, we used the same primers as the Allen Brain Atlas (<http://mouse.brain-map.org/>) (Lein et al., 2007). Standard methods for

dual color fluorescent *in situ* hybridization were used. Briefly, adult C57BL/6J brains were fresh-frozen in OCT compound and 20  $\mu$ m sections were cut. After probe hybridization and post-hybridization washes, the sections were incubated with 3% hydrogen peroxide in PBS to quench endogenous peroxidase activity. The DIG-labeled probe (candidate cell type-enriched gene) was detected with anti-DIG-POD (Roche 11207733910) followed by TSA Plus amplification (Perkin Elmer NEL745001KT). After quenching with hydrogen peroxide, the fluorescein-labeled probe (*Slc17a7*, *Pvalb*, or *Vip*) was detected with anti-Fluorescein-POD (Roche 11426346910) followed by TSA Plus amplification (Perkin Elmer NEL741001KT).

## Data analysis

Data processing steps made extensive use of Bowtie (Langmead et al., 2009; Langmead and Salzberg, 2012), Tophat (Trapnell et al., 2009), BEDTools (Quinlan and Hall, 2010), and custom scripts. All reads were aligned to the mm10 genome. Browser representations were created using AnnoJ (<http://www.annoj.org>) (Lister et al., 2009). Correlations are Pearson, unless otherwise indicated.

## RNA-seq data processing

We aligned reads from the libraries in two ways: 1) to the whole genome for genome browser visualization and 2) to the annotated transcriptome to estimate gene expression levels. RNA-seq reads were trimmed (*seqtk trimfq -b 5*) before aligning to the genome (TOPHAT v1.4.0). Gene expression levels were estimated using RSEM v1.1.20 (Li et al., 2011) calling BOWTIE v0.12.7 for protein-coding genes (mm10 iGenomes annotation). Differentially expressed genes (5% FDR) were identified through pairwise comparisons using EBSeq (v1.1) (Leng et al., 2013). Pairwise DE gene lists in Table S2 only include those genes with  $\geq 2$ -fold DE and  $\text{TPM} \geq 1$ . Additional RNA-seq data measured from fetal cortex and 6 week cortex (Lister et al., 2013) were also processed with RSEM. We used TF annotations from AnimalTFDB (Zhang et al., 2012) for all analyses focused on TFs.

## MethylC-seq data processing

MethylC-seq reads were processed as previously described using the *methylypy* pipeline (<https://bitbucket.org/schultzmatt/methylypy/> and Lister et al., 2013). Briefly, all cytosines in the forward and reverse complement strands of the mm10 reference genome (appended with the lambda phage genomic sequence) were converted to thymines followed by bowtie index building using the *build\_ref* function. The mapping of MethylC-seq reads was performed with the *run\_methylation\_pipeline* function. Adapters in MethylC-seq reads were trimmed using *cutadapt*. All cytosines in the trimmed reads were then computationally converted to thymines and mapped to a converted forward strand reference and to a converted reverse complement strand reference. Reads were only allowed to map to one location, and clonal reads were removed. The resulting datasets were stored as “alle” tables containing one row for each genomic cytosine position and columns representing the genomic context (e.g. CG, CH); the number of reads supporting a methylcytosine at that position (mc); the total number of reads at that position (h); and the calling of methylated cytosines (0 for unmethylated or 1 for methylated). The calling of methylated cytosines was performed as previously described (Lister et al., 2009; Lister et al., 2013; Schultz et al., 2015). Briefly, a binomial test was performed for each cytosine to test if the methylation level is significantly greater than 0. This was done using the bisulfite non-conversion rate derived from spiked-in unmethylated phage-lambda DNA and choosing an FDR threshold of 0.01. The bisulfite non-conversion rate was determined separately for each tri-nucleotide context. Methylated cytosines determined by the binomial test were displayed in browser tracks.



Many of our analyses are based on profiling the methylation level in CG and CH contexts (i.e. % mCG and % mCH) within a genomic region or in a set of regions. The methylation level at a set of positions  $R$  is defined as:

$$[\% mC]_R = 100 * \sum_{i \in R} mc_i / \sum_{i \in R} h_i$$

For some analyses, we adjusted these estimates to correct for bisulfite non-conversion (see below, Methods for Figure 2C-D).

### DMR finding

We estimated DMRs using a previously reported method (Ma et al., 2014; Schmitz et al., 2013; Schultz et al., 2015). Briefly, since we observed a high degree of consistency between biological replicates, we pooled reads from replicates for DMR calling to enhance the statistical power. DMR calling used five samples: E13 fetal cortex and S100b+ sorted glia (Lister et al., 2013) and pooled replicates of excitatory, PV, and VIP neurons. We used a stochastic model of our methylation data where the observed number of MethylC-seq reads supporting methylated and unmethylated cytosines at each position is taken from a binomial distribution for each sample. Since the mC base calls of unmethylated cytosines, as defined by the binomial tests, are likely to be contributed by incomplete bisulfite conversion, the mC base calls of these sites were set to zero. In order to test whether or not the observed data is consistent with the null hypothesis (i.e., the methylation level at a given position is the same in all samples), we computed a goodness-of-fit statistic (Perkins et al., 2011) using an  $N \times 2$  contingency table. This table consists of one row for each sample and a column for reads that support methylated and unmethylated cytosines, respectively. We then simulated read count data using our stochastic model under the assumptions of our null hypothesis. We used an adaptive permutation procedure to derive p-values, which ensures that sites of potential differential methylation with  $p < 0.01$  will be sampled 1,000 times; at other sites where there are 10 permutations more extreme than the original test statistic, the p-value will be greater than 0.01, and the site will not be identified as differentially methylated. We controlled the false discovery rate (FDR) at 1% by using a procedure that compares multiple sequential permutation-derived p-values (Bancroft et al., 2013). After choosing the largest p-value cutoff that satisfies  $FDR < 0.01$ , we combined significant sites (DMSs) into blocks if they had consistent methylation differences and fell within 250 bp of each other. We classified samples as hypo- or hyper-methylated based on the sign of the difference between their methylation level and the expected methylation level under the null hypothesis. Blocks that contained fewer than 4 differentially methylated sites were discarded, and the remaining blocks were called differentially methylated regions (DMRs). This method is implemented in the *DMRfind* function in *methyipy* (available at <https://bitbucket.org/schultzmatt/methyipy/>).

### DMRs and mouse strain differences

Our identification of DMRs does not factor in SNPs or indels across mouse strains. The presence of strain-specific genetic variants could potentially affect our estimates of methylation levels from MethylC-seq data. This could affect the identification of DMRs as the INTACT mice used in this study are from different genetic backgrounds. In spite of this, genome-wide DNA methylation data from excitatory neurons highly correlated with NeuN+ data from inbred C57BL/6J mice. Furthermore, we saw a high correlation between the excitatory neuron methylome and the NeuN+ methylome at localized differentially methylated regions.

To address how SNPs and indels could affect the identification of differentially methylated regions, we obtained SNPs and indels (relative to the reference C57BL/6J genome) for three mouse strains whose genomes are available (129S1, 129P2, C57BL/6N; <http://www.sanger.ac.uk/resources/mouse/genomes/>).

*R26-CAG-LSL-Sun1-sfGFP-myc* mice are on a mixed 129;C57BL/6J genetic background which includes 129S1. *PV-Cre* mice are on a mixed 129P2;C57BL/6J background. *Camk2a-Cre* mice are on a mixed C57BL/6J;C57BL/6N background. Although we do not have a similar SNP or indel list for the 129S4Sv/Jae strain (present in *VIP-Cre*) or 129X1 (also present in *R26-CAG-LSL-Sun1-sfGFP-myc*), we expect that the majority will be present in one of the other 129 strains (129S1 and 129P2). In support of this, 72.2% of all the indels and 87.5% of the SNPs that appear in either 129S1 or 129P2 strain are common to both of them (867,258/1,200,566 of indels and 4,988,081/5,697,417 SNPs).

First, by plotting the density of SNPs and indels relative to DMR locations, we found a small depletion of strain-specific variants around DMRs (data not shown), suggesting that our DMR calling is not inflated by the presence of strain-specific variants. Next, we examined 419,626 SNPs and indels that overlapped a CG position, as the DMR finder only evaluates CG sites. We then re-ran the DMR caller after removing the overlapping CG sites. 245,383 DMRs were identified using this masked data. Of these, 99.99% (245,354) overlap with the original 251,301 DMRs. Out of the original 251,301 DMRs, 97.6% (245,266) overlap with the masked DMRs. Based on this analysis, we used the original set of 251,301 DMRs in the manuscript. In Table S3, an asterisk indicates the DMRs that did not appear in the SNP-masked DMR list.

It remains possible that some strain-specific genetic variants may directly affect methylation levels. In those cases, differential methylation could be driven by strain genotype differences rather than cell type differences. Although we cannot rule this out, we think the vast majority of DMRs are cell type-driven rather than strain-driven, for the following reasons: (1) INTACT-purified excitatory neuron and NeuN+ methylomes are extremely consistent; (2) The animals are all mixed backgrounds, so strain-derived genetic components segregate independently. The high correlation between replicates argues against substantial variant effects contributing to DMR calling; and (3) The consistency of our results with known cell type markers.

### **Comparison of DMR finding using NeuN+ versus INTACT-purified excitatory neuronal nuclei**

To compare the number of total and cell type-specific DMRs identified using INTACT-purified excitatory neurons versus NeuN+ neurons from Lister et al., 2013, the MethylC-seq data for NeuN+ nuclei from the 7 week-old male (SRX314951) and the 12 month-old female (SRX314955) were combined to best match the coverage of the excitatory neuron methylomes. Identical DMR calling procedures were performed except that the NeuN+ sample was substituted for the excitatory neuron sample.

### **Identifications of UMRs, LMRs, and large DNA methylation features**

UMRs and LMRs were identified using MethylSeekR (Burger et al., 2013) with  $m = 0.5$  and 5% FDR. MethylSeekR did not classify any regions as partially methylated domains. DMVs were identified as UMRs  $\geq 5$  kb with mean.meth (column 7 in MethylSeekR output)  $\leq 15$ . To identify large hypo-DMRs, all hypo-DMRs for each cell type with inter-DMR distances less than 1 kb were merged (*bedtools merge -d 1000*). Merged hypo-DMRs were further stratified into those  $\geq 2$  kb (called “large hypo-DMRs”) and those  $< 2$  kb. mm10 CpG island annotations were downloaded from the UCSC table browser.

### **Estimation of hmC at DMVs**

To estimate the contribution of hmC to the excitatory neuron hyper-methylation of DMVs associated with *Neurog2* and *Pax6*, we mapped 6 wk cortex TAB-seq data from Lister et al., 2013 to mm10, calculated the % hmC in the region defined by Figure 7E, performed correction for non-conversion and protection, and compared it with the MethylC-seq signal of excitatory neurons in the same region.

### **ATAC-seq data processing**

Adapter sequences were trimmed from ATAC-seq reads (*cutadapt* v1.3 -a CTGTCTCTTATACACATCT -q 30 --minimum-length 36 --paired-output), before aligning (BOWTIE2 v2.1.0 -t -X2000 --no-mixed --no-discordant) and removing redundant reads (*picard MarkDuplicates*). Fragment ends were offset by 4nt towards the center of each fragment.

Peaks were called with HOMER (*findPeaks* -region -size 500 -minDist 50 -o auto -tbp 0) (Heinz et al., 2010) using sub-nucleosomal (<100 bp) fragments, and overlapping peaks were merged (*bedtools merge*). Peaks called from replicates were merged (*bedtools merge*) to yield a peak set for each cell type. We used *bedtools multiinter* to classify peaks as cell type-specific or shared, keeping only those regions greater than 100 bp. For analyses of cell type-specific versus shared peaks, the peaks following *bedtools multiinter* were used. Both the merged replicate peaks and the peaks following *bedtools multiinter* are listed in Table S3.

Footprinted sites were predicted using CENTIPEDE on ATAC-seq fragments of all lengths (Pique-Regi et al., 2011). TF binding matrices were obtained from the MEME motif database (v11, 2014 Jan 23. motif sets chen2008, hallikas2006, homeodomain, JASPAR\_CORE\_2014 Vertebrates, jolma2010, jolma2013, macisaac\_theme.v1, uniprobe\_mouse, wei2010\_mouse\_mw, wei2010\_mouse\_pbm, zhao2011) and scanned across the mouse genome to identify hits using FIMO (--output-pthresh 1E-5 --max-stored-scores 500000) (Bailey et al., 2009; Grant et al., 2010). For every ATAC-seq sample, we counted the frequency of Tn5 insertion events in 200 bp windows centered at every motif instance in the genome using *bwtool* (Pohl and Beato, 2014). These count matrices were then used by CENTIPEDE along with conservation levels at corresponding positions (phyloP score from the placental subset of the UCSC 60-way genome alignment; Karolchik et al., 2014) to learn motif-specific models of Tn5 insertion density and predict the likelihood that each motif instance across the genome is bound. We used sites predicted with greater than 95% posterior probability to be occupied as our footprint set.

To predict nucleosome positions, we employed the same procedure as Buenrostro et al., 2013. First, an estimated set of mononucleosomal fragments was generated by classifying fragments into sub-, mono-, di-, tri-, tetra-, and penta-nucleosomal fragments using a mixture of gaussians fitted to the length distribution from each sample (*mixtools* package in R). Multi-nucleosomal fragments were split into single nucleosomes by fragmenting them uniformly into the number of nucleosomes they were predicted to span, only considering those fragments whose numbers of nucleosomes were predicted with greater than 90% posterior probability. We then estimated nucleosome positioning by subtracting the “background” signal of sub-nucleosomal fragments from the “foreground” of mono-nucleosomal fragments (DANPOS -x 1 -k 1 -p 1 -a 1 -d 20 --clonalcut 0) (Chen et al., 2013).

### ChIP-seq data processing

Excitatory neuron histone modification ChIP-seq and input reads were aligned (BOWTIE v0.12.7 -m 1), and redundant reads were removed (*samtools rmdup*). CREB, SRF, and NPAS4 ChIP-seq and input reads (Kim et al., 2010) generated by SOLiD™ sequencing were aligned using BOWTIE v1.0.0 in colorspace mode (-C). FOS, FOSB, and JUNB ChIP-seq and input reads (Malik et al., 2014) were aligned using BOWTIE2 2.0.2 followed by the removal of reads with mapping quality score below 20.

We used SICER (Zang et al., 2009) to identify H3K4me1, H3K4me3, H3K27ac, and H3K27me3 ChIP-seq peaks from excitatory neurons. SICER (SICER\_V1.1) parameters with FDR = 0.001 were: redundancy threshold=1; fragment size=150; W=200, G=200 for H3K4me1, H3K4me3, and H3K27ac; and W=200, G=1000 for H3K27me3.

TF ChIP-seq peaks were identified using MACS 1.4 (Zhang et al., 2008) with a p-value cutoff at 1E-10.

### Identification of putative enhancers in excitatory neurons

Putative enhancers were defined by combining H3K27ac and H3K4me1 SICER peaks. Regions overlapping with H3K4me3 peaks or  $\pm 2.5$  kb from an annotated TSS were removed to exclude promoter features.

### **DNaseI-seq data processing**

We obtained DNaseI-seq data from 53 samples across a diverse set of neuronal and non-neuronal tissues from the mouse ENCODE project (Stamatoyannopoulos et al., 2012). Datasets were processed using the ATAC-seq analysis pipeline modified for single-end reads. Only uniquely aligning reads were kept (BOWTIE v0.12.7, options -m 1). We identified peaks of DNaseI sensitivity using HOMER (*findPeaks* -region -size 500 -mindist 50 -o auto -tbp 0). To compare DNaseI-seq peaks from whole cerebrum (GSM1014168) versus ATAC-seq peaks, we calculated the percentage of ATAC-seq peaks overlapping DNaseI-seq peaks (union of peaks in three replicates).

### **Figure-specific data analysis:**

#### **Figure S2B (comparison of INTACT and manually sorted RNA in PV neurons)**

Microarray data from manually sorted GFP+ neurons in P40 G42 transgenic mice (downloaded from GSE17806; Okaty et al., 2009) was processed using *rma* normalization from the R package *affy*. To aid in visualization, *rma*-normalized microarray values were transformed back to the linear scale, before plotting both RNA-seq TPM values and microarray values on a log scale. The Spearman correlation coefficient was calculated in R (*cor* with the option “spearman”).

#### **Figures 2C-D (global DNA methylation level)**

We adjusted the methylation level for the effect of bisulfite non-conversion, which was calibrated in each experiment by sequencing of spiked-in unmethylated phage-lambda DNA. The non-conversion rate, *s*, ranged from 0.29% to 0.38% across our samples (Table S1). We used the maximum likelihood estimate for the true methylation level (% mC) by adjusting for non-conversion as follows:

$$[\% mC]_{max.likelihood} = 100 * G\left[\frac{mc/h - s}{1 - s}\right]$$

where *mc*, *h* are the number of methylated cytosine base calls and the total cytosine base calls within a region, respectively, and  $G[x] \equiv \min[\max[x, 0], 1]$  ensures that the estimated methylation level is in the interval [0,1].

For Figure 2C, *mc* = the total number of methylated CG or CH base calls across all autosomes and *h* = the total number of all CG or CH base calls across all autosomes. For Figure 2D, the composition of methylcytosines was calculated by weighting % mCH (from Figure 2C) by the total number of CG and CH positions on autosomes.

#### **Figures 2E and S2E (line plots of % mCG and % mCH in highly cell type-specific genes)**

To assess the pattern of methylation around specific groups of DE genes, we first pooled methylation data from biological replicates. For each gene, we profiled % mCG and % mCH within 1 kb bins between 100 kb upstream of the transcription start site (TSS) and the TSS and between the transcription end site (TES) and 100 kb downstream. We divided each gene body into 10 equally spaced bins. When multiple transcripts shared the same TSS and TES, we only used 1 instance for the analysis. We then computed the

% mC (corrected for bisulfite non-conversion) for each bin. Gene lists were filtered by requiring >5 fold-change and  $\geq 0.95$  posterior probability of differential expression (PPDE) from EBSeq. We focused on genes that are differentially over-expressed in one cell type relative to both of the other cell types. These DE genes were also used for Figures 5E-F.

### **Figure S2I (comparison of intragenic and genomic % mCH across replicates and cell types)**

For each 5 kb genomic bin or gene body, we computed the CH methylation levels in each sample and corrected for bisulfite non-conversion. We excluded any genomic bins or genes with low coverage (<50 base calls) or genes with short (<500 bp TSS-TES) length. We normalized the mCH level for each sample by the median across all genomic bins or gene bodies. We then computed the ratio of the methylation levels between cell types (solid lines). As a control, we also computed this ratio for comparisons of biological replicates (dashed lines). In comparisons where both samples had very low levels of mCH (<0.5%), we set the fold-change to 1.

### **Figures 3B, D (Venn diagrams of ATAC-seq peaks and DMRs)**

Venn diagrams were created using *eulerAPE* (Micallef and Rodgers, 2014).

### **Figure S3B, middle and right (activity-dependent TF peaks at hypo-DMRs and differential ATAC-seq peaks)**

For TF  $i$ , the enrichment or depletion of each hypo-DMR category  $j$  overlapping each TF  $i$  ChIP-seq peak category, relative to all DMRs, was represented as:  $\log_2(\text{fraction of category } j \text{ hypo-DMRs overlapping TF } i \text{ peaks} / \text{fraction of all DMRs overlapping TF } i \text{ peaks})$ . The hypergeometric test (MATLAB *hygecdf*) was used to test for significance: *hygecdf(number of category j hypo-DMRs associated with TF i ChIP-seq peaks, sample size of all DMRs, number of all DMRs associated with TF i ChIP-seq peaks, sample size of category j hypo-DMRs)*. The option 'upper' was applied for testing enrichment. An analogous test was used to assess the enrichment or depletion of each cell type-specific ATAC-seq peak category overlapping each TF ChIP-seq peak.

### **Figures 3E and S3C (levels of CG and CH DNA methylation, ATAC-seq reads, histone ChIP-seq reads, and activity-dependent TF ChIP-seq reads at DMRs)**

Hypo-methylated cell types in DMRs were identified from *methyipy*, and only DMRs with one or two hypo-methylated cell types were displayed in the heatmap. Wiggle files for DNA methylation levels were created from *methyipy run\_methylation\_pipeline* output files at 100 bp resolution for CG and CH contexts. Wiggle files for ATAC-seq were generated by the pileup of sub-nucleosomal (<100 bp) reads at 100 bp resolution. Wiggle files for histone modification and TF ChIP-seq were generated using MACS14 with options -w and --space 100. Profiles of DNA methylation, ATAC-seq (normalized for library size), and ChIP-seq (normalized for library size) were plotted in a 3 kb region centered at DMRs using the wiggle files.

### **Figure S3D (correlation between CG and CH methylation at DMRs)**

% mCG and % mCH levels for excitatory, PV, and VIP neurons in each DMR were normalized by the mean % mCG and % mCH level across the three cell types for that DMR. The Pearson correlation between normalized mCG and mCH was calculated with the MATLAB *corr* function.

### **Figure 4B (correlation across cell types of ATAC-seq density at peaks)**

The similarity of ATAC-seq read distributions between pairs of ATAC-seq samples was quantified using the Pearson correlation of read densities over the union of peaks called across all samples (*deeptools bamCorrelate*) (Ramírez et al., 2014).

### **Figure S5A (evaluation of DNA methylation at footprints)**

To investigate the correlation between TF binding and local DNA methylation, we focused on footprints that were unique to one of the three cell types; that is, we excluded footprints that had the same start and end site in more than one cell type. Then, for each footprint of a given TF, we calculated methylation levels (% mCG, % mCH, and % mCA) within  $\pm 50$  bp of the footprint start site; this value is the methylation level for “FP present” locations (y-axis). We compared this with the methylation levels at the same location in the other cell types (“FP absent,” x-axis), provided that the TF is expressed ( $\text{TPM} \geq 30$ ) in these other cells. Then, the average methylation level using all MethylC-seq reads from footprinted regions (i.e., excitatory reads at excitatory-specific FPs, PV reads at PV-specific FPs, and VIP reads at VIP-specific FPs) was plotted against the average methylation level of all MethylC-seq reads from non-footprinted regions (i.e., excitatory and PV reads at VIP-specific FPs, etc.).

### **Figures 5B and S5B (enrichment of TF footprints and hypo-DMR motifs)**

For Figure 5B, we ranked each TF by the relative enrichment of their footprints in a foreground category of cell type-specific ATAC-seq peaks versus a background of the other two categories of cell type-specific ATAC-seq peaks (e.g., footprints in excitatory-only peaks versus PV-only and VIP-only peaks) and additionally required that the TF itself be expressed ( $\geq 30$  TPM) in the foreground cell type (e.g., in excitatory neurons). The significance of enrichment was estimated using the pairwise Fisher’s test (*pairwise.Fisher.test* in the *fsmB* R package) (i.e., to compare the ratio of a TF’s footprints to the total number of footprints predicted in one cell type against the corresponding ratio computed from footprints in the other two cell types). The same test was used to compare the enrichment of a TF’s motifs in one category of cell type-specific hypo-DMRs versus the other two categories of cell type-specific hypo-DMRs.

To assess TF motifs that were enriched in DMRs hypo-methylated in one or two INTACT-purified cell types as well as fetal cortex and glia (Figure S5B), hypergeometric tests were performed for each TF motif using the occurrence of the motif in all DMRs as the background.

### **Figures 5D (construction of TF-TF regulatory networks)**

TF A was predicted to regulate TF B when: (1) TF A was expressed in a cell type ( $\geq 30$  TPM), (2) TF A had a predicted footprint (FP A) in a cell type-specific ATAC-seq peak, (3) the ATAC-seq peak was within 10 kb of the TSS for TF B, and (4) TF B was expressed in that cell type ( $\geq 30$  TPM). The resulting set of predicted regulatory interactions was visualized as a network (*igraph* package in R), omitting TFs with more than 20 connections to ease visualization. To define a pan-neuronal regulatory network, we identified footprints common to all three cell types that occurred in shared ATAC-seq peaks and did not overlap ubiquitous DNaseI peaks (peaks occurring in at least 40 out of 53 processed DNaseI-seq samples). The full networks are listed in Table S4.

### **Figures S6A-B (sparse generalized linear model of mRNA expression)**

To assess how well mRNA expression levels correlate with a combination of epigenetic and chromatin features, we fit a generalized linear model using the MATLAB implementation of *cvglmnet* (Friedman et al., 2010) with the Poisson distribution and parameter  $\alpha = 1$  (corresponding to LASSO regression). This model assumes Poisson distributed noise and uses LASSO regularization to promote sparseness, i.e. to fit the model using a small subset of features. We used 10-fold cross-validation to avoid overfitting and



default values for all other parameters. For each gene, we used the longest isoform to guarantee that each gene contributes only once to the dataset and there is no overlap between training and test sets. To define features for this analysis, we created 18 windows of varying sizes, ranging from 200 bp to 32 kb, surrounding each TSS (Table S5). For mCG and mCH we also included two additional windows for the gene body and the flanking region. Within each window we computed the value of mCG, mCH, ATAC-seq, and DMR density, resulting in a total of  $4 \times 18 + 2 \times 2 = 76$  features. Using 7 parameters in the regression model achieves 1 standard error above the minimum cross-validated error.

### **Figures 6B-E and S6C (k-means clustering of genes by intragenic mCH followed by assessment of gene expression and ATAC-seq peak enrichment in each cluster)**

To identify sets of genes that share similar DNA methylation patterns in an unbiased fashion, we applied k-means clustering to the gene body mCH. We profiled % mCH in gene bodies (TSS-TES) within each of eight samples included in this analysis (Fetal and Adult Cortex, NeuN+, NeuN-, and Glia from Lister et al., 2013; Exc, PV, and VIP from the current study). We excluded 468 genes with short gene bodies (<500 bp TSS-TES) or with low coverage in our methylome datasets (<50 cytosine base calls within the gene body in any sample); the remaining 23,023 genes were included. When multiple transcripts shared the same TSS and TES, we only used 1 instance for the analysis. To compensate for the differing genome-wide background level of % mCH in different cell types, we normalized the methylation level in each sample for each gene by the average over the gene's distal flanking region (50-100 kb upstream of TSS or downstream of TES). We then log-transformed the normalized methylation levels. Next, we used the MATLAB function *kmeans* to apply k-means clustering using data from five samples representing distinct cell types or developmental stages (Fetal cortex, Exc, PV, VIP, and Glia). Clustering used 1 minus the correlation coefficient of normalized mCH values across genes as a distance measure. We chose to extract  $k=25$  clusters to capture a diverse range of methylation features, while still allowing visualization and statistical enrichment analysis of functional association for each gene set (Figure S6C). We repeated the clustering procedure five times using random initialization of the cluster centers, choosing as the final estimate the run with the smallest within-cluster sum of distances from each point to the cluster centroid.

To display the CH methylation patterns within these gene clusters in Figure 6B, we profiled % mCH in 1 kb bins starting 100 kb upstream of the TSS and ending 100 kb downstream of the TES. To compare genes with different lengths, we divided each gene body into 10 non-overlapping bins of equal size extending from the TSS to the TES. Methylation levels were normalized by the flanking region as described above. We then linearly interpolated the gene-body mCH data at 100 evenly spaced bins within the gene body in order to give similar visual weight to the gene-body and flanking methylation data. Finally, we smoothed and downsampled the genes 40-fold to allow representation of genome-scale features.

RNA TPM levels were plotted for each cluster (Figure 6C). Last, we assessed the enrichment of specific categories of DE genes (Figure 6D) and ATAC-seq peaks located within  $\pm 10$  kb of each gene's TSS (Figure 6E) within each cluster using hypergeometric test with Benjamini-Hochberg FDR control.

### **Figures 6F and S6D, top (DNA methylation levels relative to nucleosome calls)**

For each nucleosome call, we counted the number of sequenced CG (Figure S6D, top) or CH (Figure 6F) base calls starting from the nucleosome center up to 2000 bp away. We also counted the number of these sequenced base calls that were methylated. The ratio of the methylated to the total number of sequenced base calls is the average mCG or mCH level at that position. Because we were able to average over all of the estimated nucleosome positions, binning was not necessary and the mCH level was estimated as a function of distance with 1 bp resolution. We determined the average in the flanking region by summing

over all base calls from 1-2 kb upstream and downstream of each nucleosome call. This average is a single number, not a function of distance from the nucleosome. The normalized curves in the figures show the mCG and mCH level divided by the flanking region average.

### **Figures S7C-D (GC content and CG methylation level of hypo-methylated features)**

GC contents were computed with *bedtools nuc* (Figure S7C). Genomic regions matching the sizes of excitatory hypo-DMRs < 2 kb were randomly selected with *bedtools shuffle*. The random selection was repeated 100 times.

### **Figure 7C (enrichment of histone marks, ATAC-seq reads, and RNA over hypo-methylated features)**

For excitatory neurons, we divided DMVs into those that overlapped with SICER peaks for H3K4me3 by  $\geq 1$  bp or those that overlapped with SICER peaks for H3K27me3 by  $\geq 1$  bp. For each type of DNA methylation feature in excitatory neurons, the  $\log_2$  enrichment of each histone modification over the input (both normalized for library size) was plotted (left). For ATAC-seq,  $\log_2(1 + \text{ATAC-seq} < 100 \text{ bp reads pileup per } 10 \text{ million reads})$  was plotted. Protein-coding genes were associated with all hypo-methylated features if  $\geq 1$  bp overlap was found between these features and the region from 10 kb upstream of the TSS to the TES.  $\log_2(\text{TPM} + 1)$  values were plotted for all associated genes.

### **Figure 7D (overlap of DE genes in hypo-methylated regions)**

Differentially expressed (DE) genes identified by EBSeq with  $\geq 2$  fold-change were used for this analysis. Protein-coding genes were associated with all hypo-methylated features identically to Figure 7C. The significance of the overlap between each DNA methylation feature and DE genes was tested by hypergeometric distribution using MATLAB *hygecdf* function with *hygecdf(number of category i feature overlapping with category j DE genes, sample size of all DNA methylation features, number of all DNA methylation features that overlap with category j DE genes, sample size of category i feature)*. The option 'upper' was applied for testing enrichment. All DNA methylation features were defined as combined large hypo-DMRs, merged hypo-DMRs with length less than 2 kb, and DMVs identified for Exc, PV, and VIP neurons. Significance was set at  $q < 1E-5$ .

### **Figure S7G (gene ontology enrichment in DMVs)**

GO enrichment for each group of DMVs was performed using GREAT (McLean et al., 2010). For the background, DMVs were combined with UMRs between 1-3.5 kb and mean mCG  $\leq 15\%$ .

### **Supplemental References**

Bailey, T.L., Boden, M., Buske, F.A., Frith, M., Grant, C.E., Clementi, L., Ren, J., Li, W.W., and Noble, W.S. (2009). MEME SUITE: tools for motif discovery and searching. *Nucleic Acids Res.* 37, W202-208.

Bancroft, T., Du, C., and Nettleton, D. (2013). Estimation of false discovery rate using sequential permutation p-values. *Biometrics* 69, 1–7.

Chen, K., Xi, Y., Pan, X., Li, Z., Kaestner, K., Tyler, J., Dent, S., He, X., and Li, W. (2013). DANPOS: dynamic analysis of nucleosome position and occupancy by sequencing. *Genome Res.* 23, 341-351.

- Dugas, J.C., Tai, Y.C., Speed, T.P., Ngai, J., and Barres, B.A. (2006). Functional genomic analysis of oligodendrocyte differentiation. *J. Neurosci.* *26*, 10967-10983.
- Friedman, J., Hastie, T., and Tibshirani, R. (2010). Regularization paths for generalized linear models via coordinate descent. *J. Stat. Softw.* *33*, 1-22.
- Garber, M., Yosef, N., Goren, A., Raychowdhury, R., Thielke, A., Guttman, M., Robinson, J., Minie, B., Chevrier, N., Itzhaki, Z., et al. (2012). A high-throughput chromatin immunoprecipitation approach reveals principles of dynamic gene regulation in mammals. *Mol. Cell* *47*, 810-822.
- Grant, C.E., Bailey, T.L., and Noble, W.S. (2011). FIMO: scanning for occurrences of a given motif. *Bioinformatics* *27*, 1017-1018.
- Heinz, S., Benner, C., Spann, N., Bertolino, E., Lin, Y.C., Laslo, P., Cheng, J.X., Murre, C., Singh, H., and Glass, C.K. (2010). Simple combinations of lineage-determining transcription factors prime cis-regulatory elements required for macrophage and B cell identities. *Mol. Cell* *38*, 576-589.
- Karolchik, D., Barber, G.P., Casper, J., Clawson, H., Cline, M.S., Diekhans, M., Dreszer, T.R., Fujita, P.A., Guruvadoo, L., Haeussler, M., et al. (2014). The UCSC Genome Browser database: 2014 update. *Nucleic Acids Res.* *42*, D764-770.
- Langmead, B., and Salzberg, S.L. (2012). Fast gapped-read alignment with Bowtie 2. *Nat. Methods* *9*, 357-359.
- Langmead, B., Trapnell, C., Pop, M., and Salzberg, S.L. (2009). Ultrafast and memory-efficient alignment of short DNA sequences to the human genome. *Genome Biol.* *10*, R25.
- Lein, E.S., Hawrylycz, M.J., Ao, N., Ayres, M., Bensinger, A., Bernard, A., Boe, A.F., Boguski, M.S., Brockway, K.S., Byrnes, E.J., et al. (2007). Genome-wide atlas of gene expression in the adult mouse brain. *Nature* *445*, 168-176.
- Leng, N., Dawson, J.A., Thomson, J.A., Ruotti, V., Rissman, A.I., Smits, B.M., Haag, J.D., Gould, M.N., Stewart, R.M., and Kendzioriski, C. (2013). EBSeq: an empirical Bayes hierarchical model for inference in RNA-seq experiments. *Bioinformatics* *29*, 1035-1043.
- Li, B., and Dewey, C.N. (2011). RSEM: accurate transcript quantification from RNA-Seq data with or without a reference genome. *BMC Bioinformatics* *12*, 323.
- Lister, R., Pelizzola, M., Dowen, R.H., Hawkins, R.D., Hon, G., Tonti-Filippini, J., Nery, J.R., Lee, L., Ye, Z., Ngo, Q.M., et al. (2009). Human DNA methylomes at base resolution show widespread epigenomic differences. *Nature* *462*, 315-322.
- Ma, H., Morey, R., O'Neil, R.C., He, Y., Daughtry, B., Schultz, M.D., Hariharan, M., Nery, J.R., Castanon, R., Sabatini, K., et al. (2014). Abnormalities in human pluripotent cells due to reprogramming mechanisms. *Nature* *511*, 177-183.
- McLean, C.Y., Bristor, D., Hiller, M., Clarke, S.L., Schaar, B.T., Lowe, C.B., Wenger, A.M., and Bejerano, G. (2010). GREAT improves functional interpretation of cis-regulatory regions. *Nat. Biotechnol.* *28*, 495-501.
- Micallef, L., and Rodgers, P. (2014). eulerAPE: drawing area-proportional 3-Venn diagrams using ellipses. *PLoS One* *9*, e101717.

- Okaty, B.W., Miller, M.N., Sugino, K., Hempel, C.M., and Nelson, S.B. (2009). Transcriptional and electrophysiological maturation of neocortical fast-spiking GABAergic interneurons. *J. Neurosci.* *29*, 7040-7052.
- Pédélecq, J.D., Cabantous, S., Tran, T., Terwilliger, T.C., and Waldo, G.S. (2006). Engineering and characterization of a superfolder green fluorescent protein. *Nat. Biotechnol.* *24*, 79-88.
- Perkins, W., Tygert, M., and Ward, R. (2011). Computing the confidence levels for a root-mean-square test of goodness-of-fit. *App. Math and Comp.* *217*, 9072-9084.
- Pique-Regi, R., Degner, J.F., Pai, A.A., Gaffney, D.J., Gilad, Y., and Pritchard, J.K. (2011). Accurate inference of transcription factor binding from DNA sequence and chromatin accessibility data. *Genome Res.* *21*, 447-455.
- Pohl, A., and Beato, M. (2014). bwtool: a tool for bigWig files. *Bioinformatics* *30*, 1618-1619.
- Quinlan, A.R., and Hall, I.M. (2010). BEDTools: a flexible suite of utilities for comparing genomic features. *Bioinformatics* *26*, 841-842.
- Ramírez, F., Dündar, F., Diehl, S., Grüning, B.A., and Manke, T. (2014). deepTools: a flexible platform for exploring deep-sequencing data. *Nucleic Acids Res.* *42*, W187-191.
- Schmitz, R.J., Schultz, M.D., Urich, M.A., Nery, J.R., Pelizzola, M., Libiger, O., Alix, A., McCosh, R.B., Chen, H., Schork, N.J., et al. (2013). Patterns of population epigenomic diversity. *Nature* *495*, 193-198.
- Schultz, M.D., He, Y., Whitaker, J.W., Hariharan, M., Mukamel, E.A., Leung, D., Rajagopal, N., Nery, J.R., Urich, M.A., Chen, H., et al. (2015). Human body epigenome maps reveal noncanonical DNA methylation variation. *Nature* doi: 10.1038/nature14465.
- Soriano, P. (1999). Generalized lacZ expression within the ROSA26 Cre reporter strain. *Nat. Genet.* *21*, 70-71.
- Trapnell, C., Pachter, L., and Salzberg, S.L. (2009). TopHat: discovering splice junctions with RNA-Seq. *Bioinformatics* *25*, 1105-1111.
- Zang, C., Schones, D.E., Zeng, C., Cui, K., Zhao, K., and Peng, W. (2009). A clustering approach for identification of enriched domains from histone modification ChIP-Seq data. *Bioinformatics* *25*, 1952-1958.
- Zhang, H.M., Chen, H., Liu, W., Liu, H., Gong, J., Wang, H., and Guo, A.Y. (2012). AnimalTFDB: a comprehensive animal transcription factor database. *Nucleic Acids Res.* *40*, D144-149.
- Zhang, Y., Liu, T., Meyer, C.A., Eeckhoute, J., Johnson, D.S., Bernstein, B.E., Nusbaum, C., Myers, R.M., Brown, M., Li, W., et al. (2008). Model-based analysis of ChIP-Seq (MACS). *Genome Biol.* *9*, R137.